

# Analysis Methods in Atmospheric and Oceanic Science

AOSC 652

Least Squares Analysis, Statistical Regression,  
and Spline Fitting:  
Day 1

26 Sep 2016

# AOSC 652: Analysis Methods in AOSC

## Assignment 4a: Integer Sorting

Only 9 or 15 students turned in code that was working: output to file fort.99 designed to let you check if your code worked

If you know your code is not working, please ask me, Jeff, or Walt for help

We had hoped that the answers to the functional dependence of the timing would be drawn from the readings: please try to do the readings

# AOSC 652: Analysis Methods in AOSC

Review Mon:

Call to piksrt in our code:

```
call piksrt(iarray_in,iarray_out,npts)
```

New subroutine piksrt to comply with our call statement

```
subroutine piksrt(arr_in,arr_out,n)
integer n,i,j
real a,arr_in(n),arr_out(n)
do i=1,n
  arr_out(i)=arr_in(i)
enddo
do j=2,n          ! Pick out each element in turn
  a=arr_out(j)
  do i=j-1,1,-1  ! Look for the place to insert it
    if(arr_out(i).le.a) goto 10
    arr_out(i+1)=arr_out(i)
  enddo
  i=0
10  arr_out(i+1)=a      ! Insert it
enddo
return
end
```

# AOSC 652: Analysis Methods in AOSC

Review Mon:

Call to piksrt in our code:

```
call piksrt(iarray_in,iarray_out,npts)
```

New subroutine piksrt to comply with our call statement

```
subroutine piksrt(arr_in,arr_out,n)
integer n,i,j
integer a,arr_in(n),arr_out(n)
do i=1,n
  arr_out(i)=arr_in(i)
enddo
do j=2,n          ! Pick out each element in turn
  a=arr_out(j)
  do i=j-1,1,-1  ! Look for the place to insert it
    if(arr_out(i).le.a) goto 10
    arr_out(i+1)=arr_out(i)
  enddo
  i=0
10  arr_out(i+1)=a      ! Insert it
enddo
return
end
```

# AOSC 652: Analysis Methods in AOSC

## Assignment 4a: Integer Sorting

Only 9 or 15 students turned in code that was working: output to file fort.99 designed to let you check if your code worked

If you know your code is not working, please ask me, Jeff, or Walt for help

We had hoped that the answers to the functional dependence of the timing would be drawn from the readings: please try to do the readings

Reading: piksrt speed varies as  $N^2 \Rightarrow$  sort of 1 million points would take  $10 \times 10 \times$  speed of 100,000 point sort or  $\sim 100$  min

# AOSC 652: Analysis Methods in AOSC

## Assignment 4a: Integer Sorting

Only 9 or 15 students turned in code that was working: output to file fort.99 designed to let you check if your code worked

If you know your code is not working, please ask me, Jeff, or Walt for help

We had hoped that the answers to the functional dependence of the timing would be drawn from the readings: please try to do the readings

Reading: piksrt speed varies as  $N^2 \Rightarrow$  sort of 1 million points would take  $10 \times 10 \times$  speed of 100,000 point sort or  $\sim 100$  min

heapsort speed varies as  $N \log_2 N$

# AOSC 652: Analysis Methods in AOSC

## Assignment 4a: Integer Sorting

Only 9 or 15 students turned in code that was working: output to file fort.99 designed to let you check if your code worked

If you know your code is not working, please ask me, Jeff, or Walt for help

We had hoped that the answers to the functional dependence of the timing would be drawn from the readings: please try to do the readings

Reading: piksrt speed varies as  $N^2 \Rightarrow$  sort of 1 million points would take  $10 \times 10 \times$  speed of 100,000 point sort or  $\sim 100$  min

heapsort speed varies as  $N \log_2 N$

Ratio of run times either:

$$N^2 / N \log_2 N = 10^6 / \log_2 (10^6) = \sim 50,000 \text{ (theory)}$$

\*or\*

$$100 \text{ min} \times 60 \text{ sec} / \text{min} / 2.6 \text{ sec} = \sim 2300 \text{ (observation)}$$

# AOSC 652: Analysis Methods in AOSC

## Assignment 4a: Integer Sorting

It is remarkable that, for a little bit more effort (~30 min for most students, hopefully), we can obtain an algorithm thousands of times more efficient for sorting a million numbers than a “brute force” method

There are many times in your research life you may need to make a choice between “clock time” in getting a computation started and “run time” on the computer

Reading: piksrt speed varies as  $N^2 \Rightarrow$  sort of 1 million points would take  
 $10 \times 10 \times$  speed of 100,000 point sort or ~100 min

heapsort speed varies as  $N \log_2 N$

Ratio of run times either:

$$N^2 / N \log_2 N = 10^6 / \log_2 (10^6) = \sim 50,000 \text{ (theory)}$$

\*or\*

$$100 \text{ min} \times 60 \text{ sec} / \text{min} / 2.6 \text{ sec} = \sim 2300 \text{ (observation)}$$



# AOSC 652: Analysis Methods in AOSC

Please take Assignment 4B, check your numbers, and return on Wed at start of class

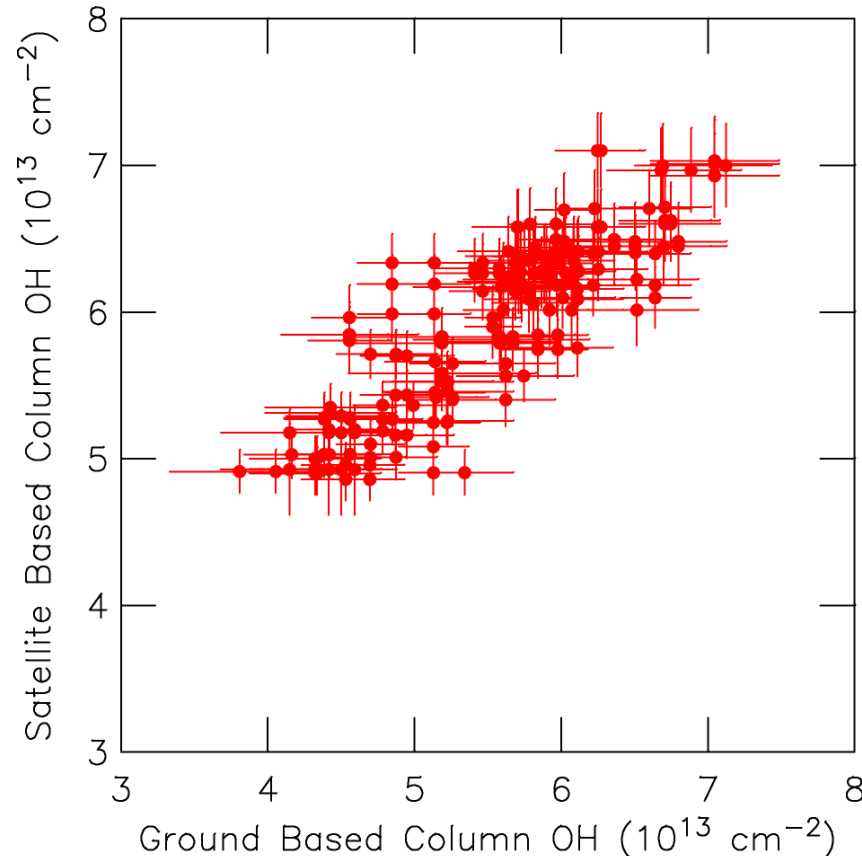
Name: \_\_\_\_\_

Year	Temperature
1951	
1952	
1953	
1954	
1955	
1956	
1957	
1958	
1959	
1960	
1961	
1962	
1963	
1964	
1965	
1966	

Year	Temperature
1967	
1968	
1969	
1970	
1971	
1972	
1973	
1974	
1975	
1976	
1977	
1978	
1979	
1980	
Baseline Hand	
Baseline Code	

# AOSC 652: Analysis Methods in AOSC

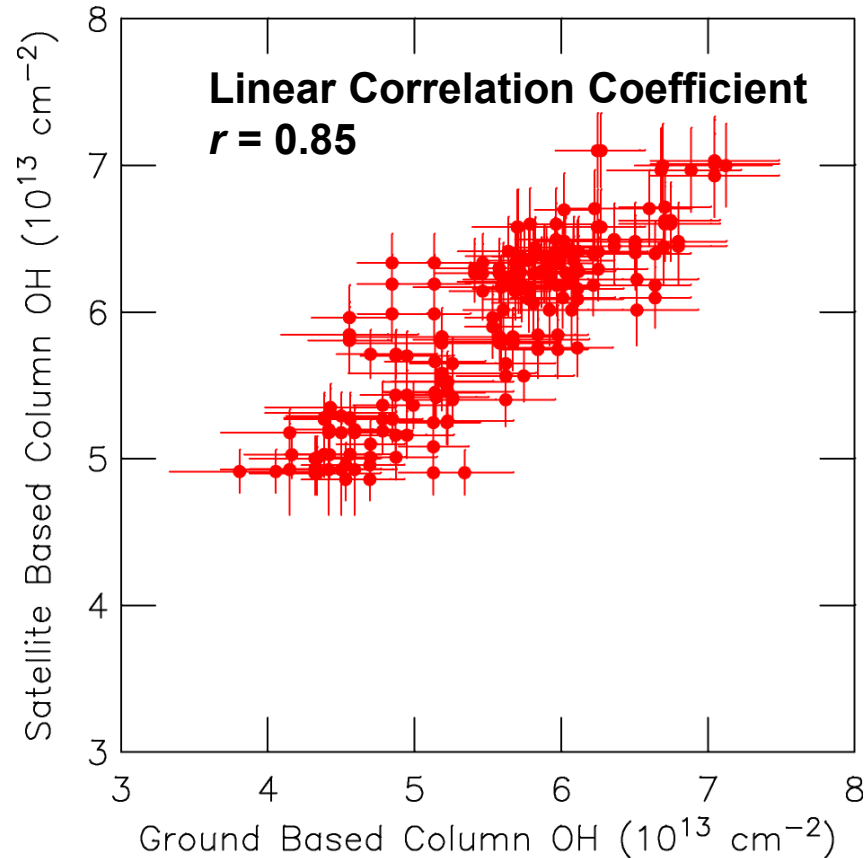
Suppose you have two sets of measurements (or data and model) that you'd like to relate.



**What are some aspects of the data that are typically examined?**

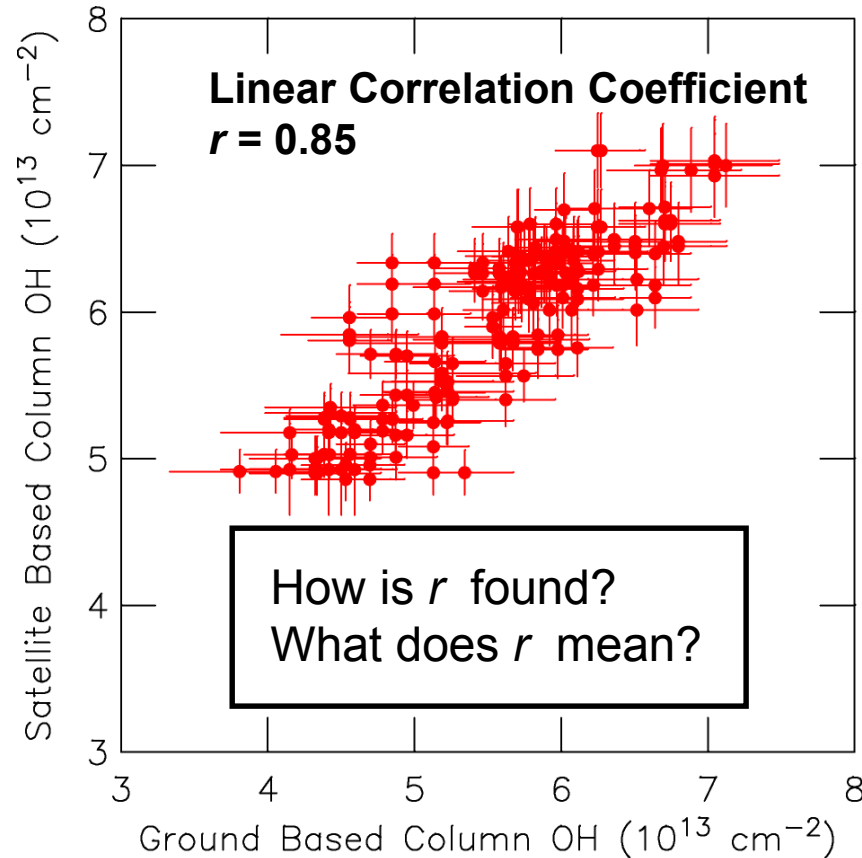
# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



# AOSC 652: Analysis Methods in AOSC

## Linear Correlation Coefficient:

$$r = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}}$$

$r$  must lie between  $-1$  and  $1$

If  $r = 1$ , the data are said to have a *complete positive correlation*

$r = 0$ , the data are said to be *uncorrelated*

# AOSC 652: Analysis Methods in AOSC

## Linear Correlation Coefficient:

$$r = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}}$$

$r$  must lie between  $-1$  and  $1$

If  $r = 1$ , the data are said to have a *complete positive correlation*

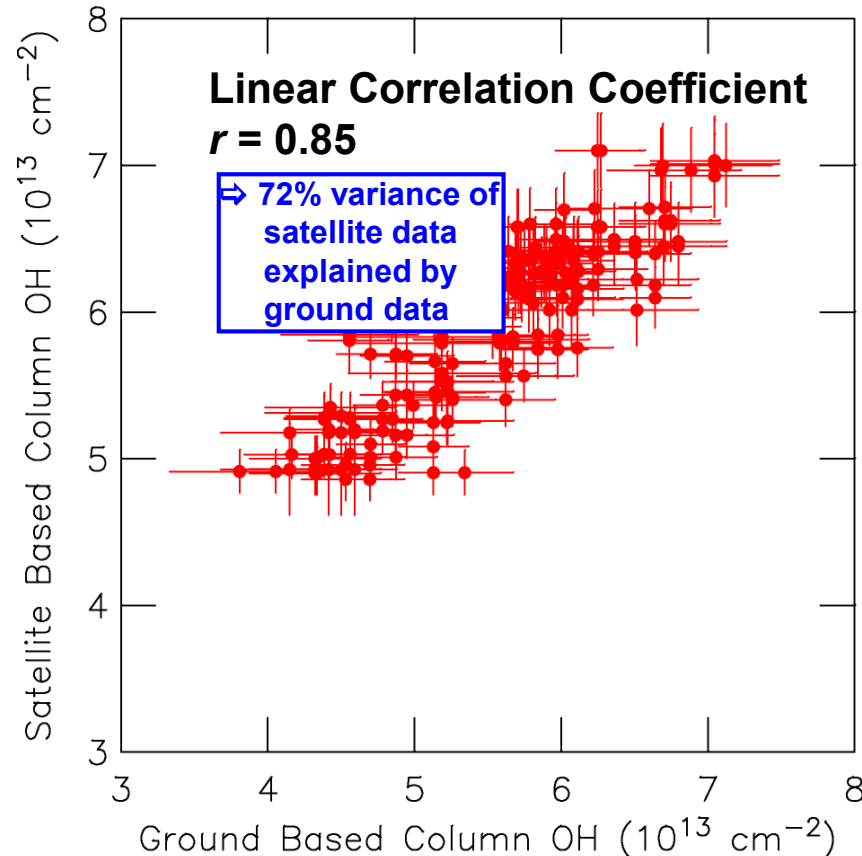
$r = 0$ , the data are said to be *uncorrelated*

$r^2 \times 100 =$  percent of variance in common between  $x$  and  $y$

See <http://www.mega.nu/ampp/rummel/uc.htm#C2> for a nice explanation

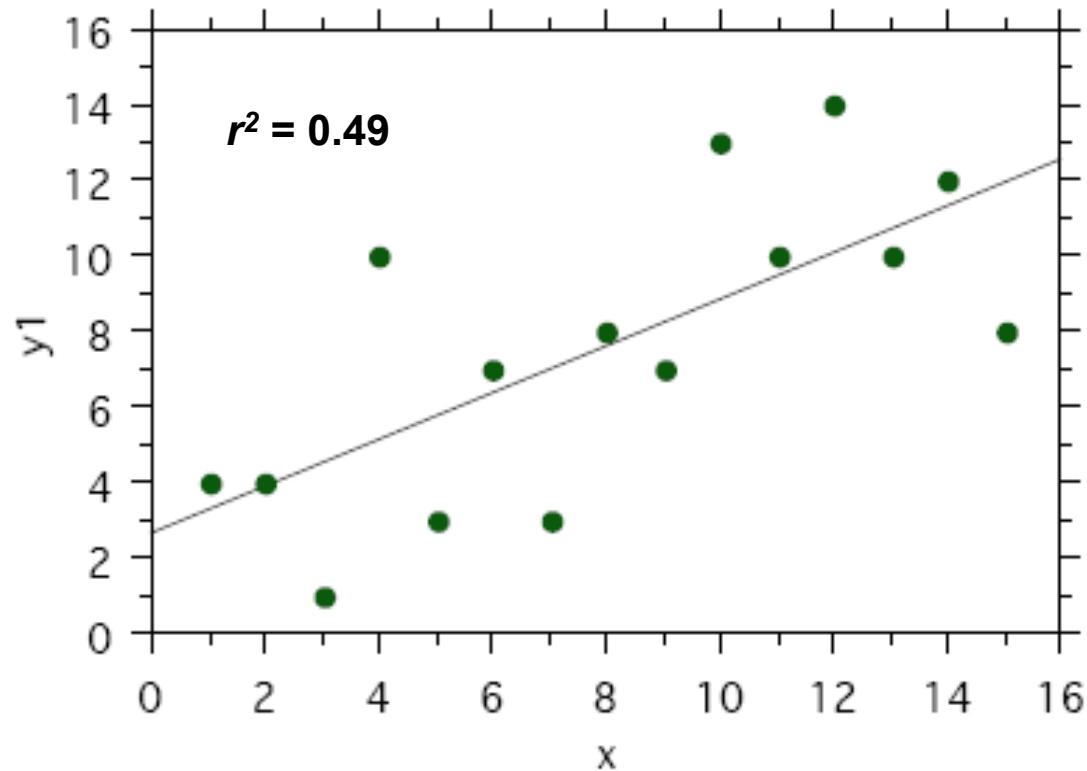
# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



# AOSC 652: Analysis Methods in AOSC

## Visualization of data is vital



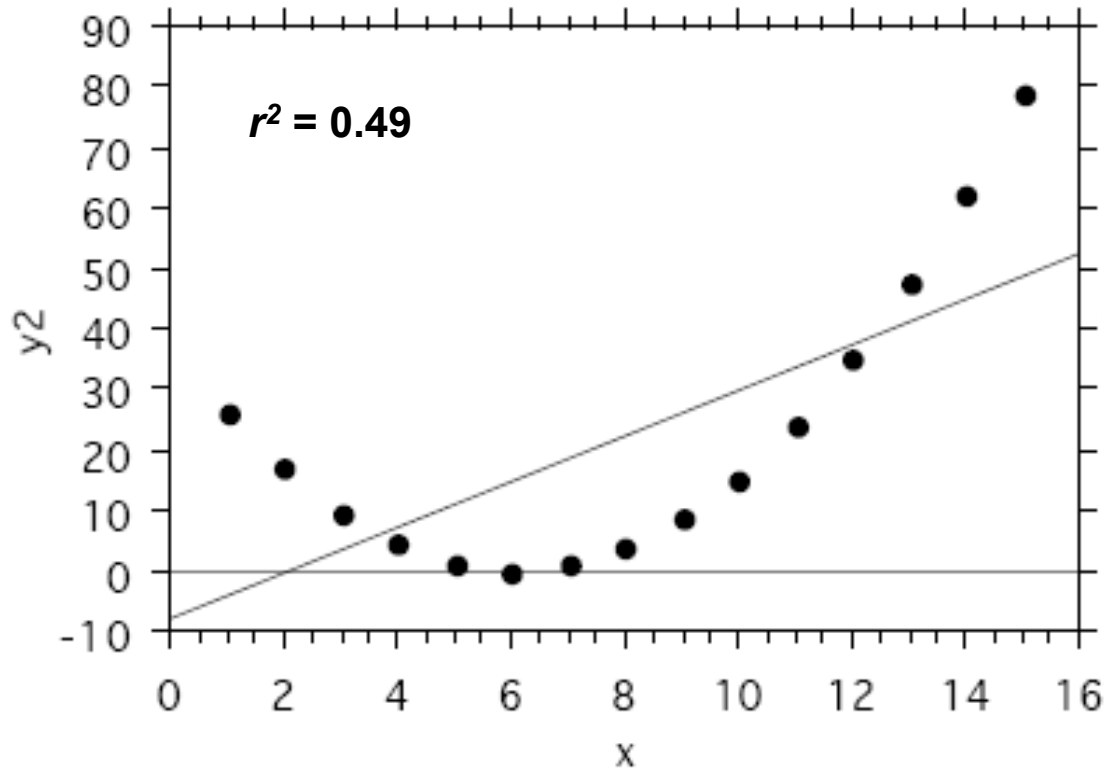
From Dennis Hartmann's class notes:

[http://www.atmos.washington.edu/~dennis/552\\_Notes\\_3.pdf](http://www.atmos.washington.edu/~dennis/552_Notes_3.pdf)



# AOSC 652: Analysis Methods in AOSC

## Visualization of data is vital

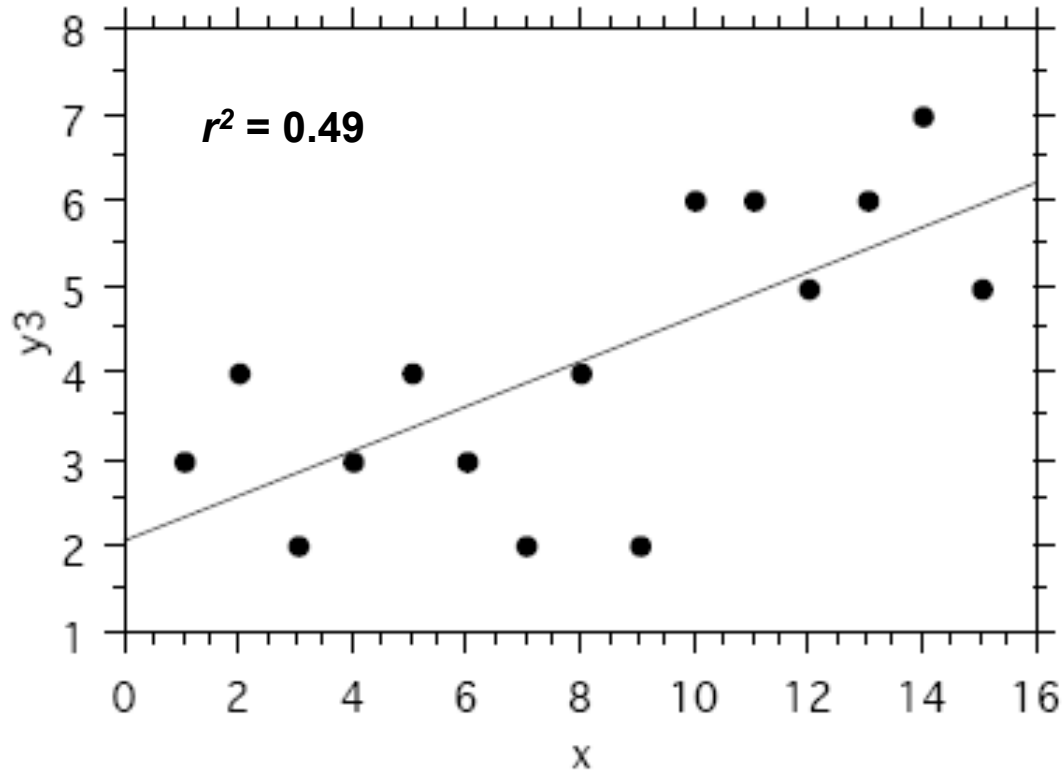


From Dennis Hartmann's class notes:

[http://www.atmos.washington.edu/~dennis/552\\_Notes\\_3.pdf](http://www.atmos.washington.edu/~dennis/552_Notes_3.pdf)

# AOSC 652: Analysis Methods in AOSC

## Visualization of data is vital

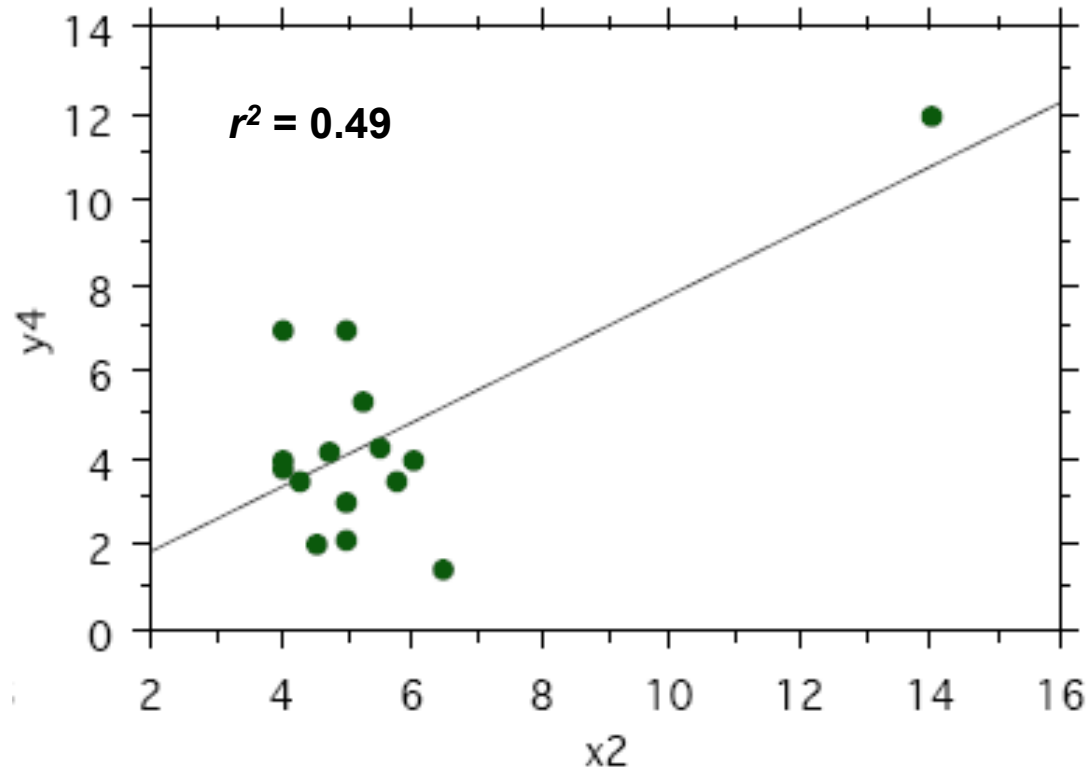


From Dennis Hartmann's class notes:

[http://www.atmos.washington.edu/~dennis/552\\_Notes\\_3.pdf](http://www.atmos.washington.edu/~dennis/552_Notes_3.pdf)

# AOSC 652: Analysis Methods in AOSC

## Visualization of data is vital

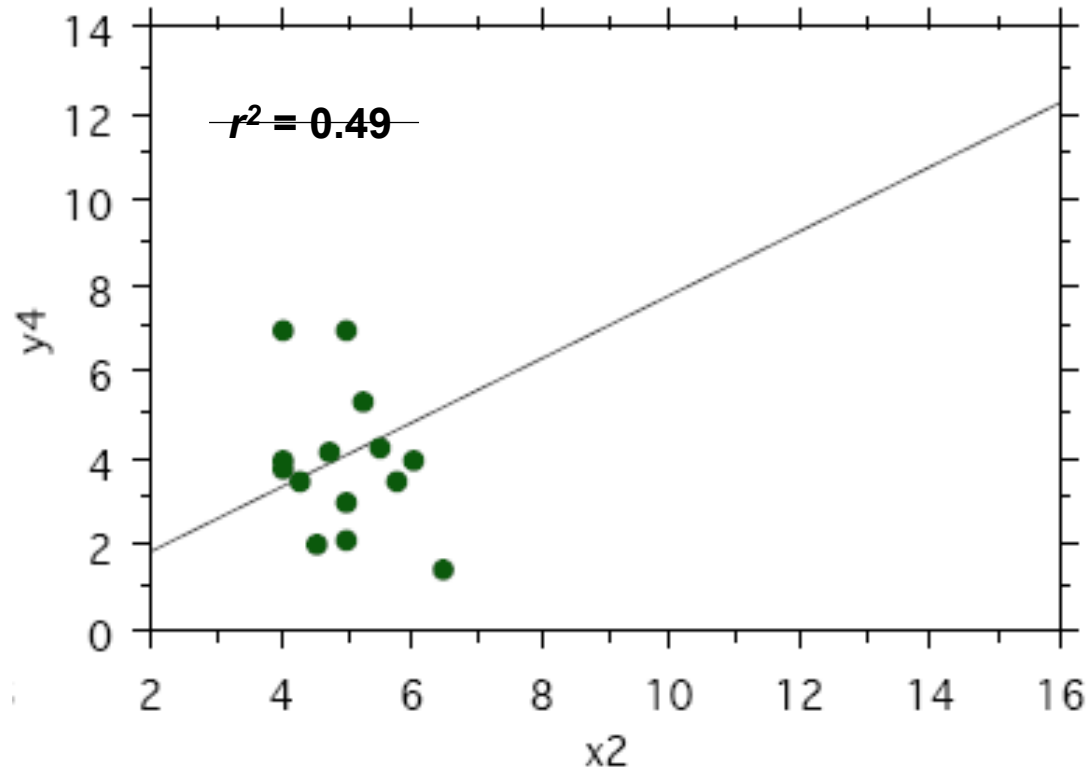


From Dennis Hartmann's class notes:

[http://www.atmos.washington.edu/~dennis/552\\_Notes\\_3.pdf](http://www.atmos.washington.edu/~dennis/552_Notes_3.pdf)

# AOSC 652: Analysis Methods in AOSC

## Visualization of data is vital

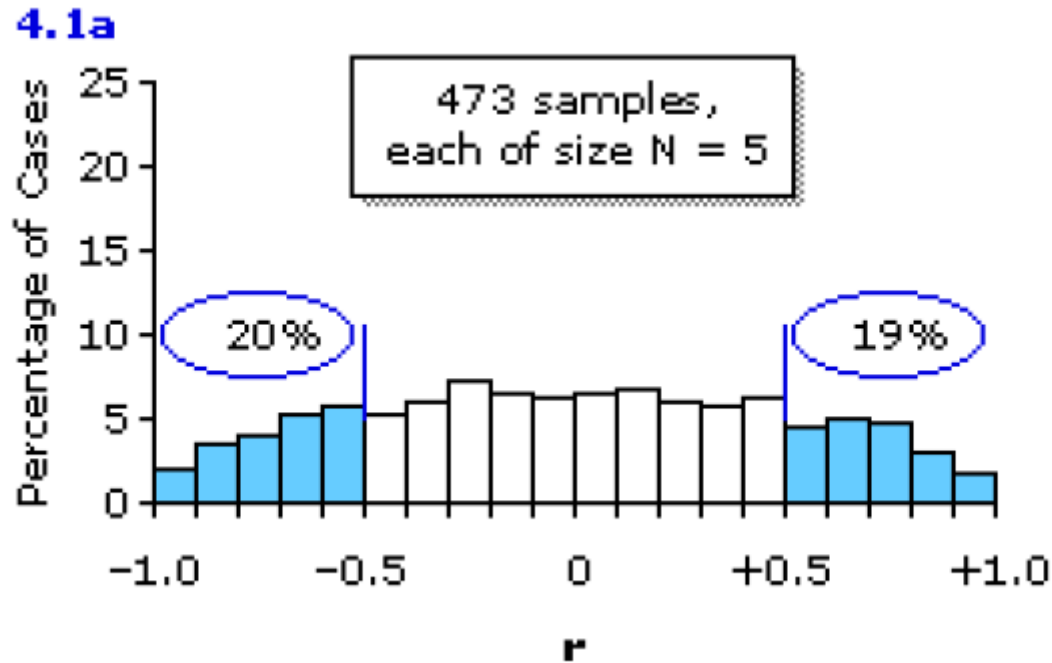


From Dennis Hartmann's class notes:

[http://www.atmos.washington.edu/~dennis/552\\_Notes\\_3.pdf](http://www.atmos.washington.edu/~dennis/552_Notes_3.pdf)

# AOSC 652: Analysis Methods in AOSC

## Test for statistical significance of correlation vital



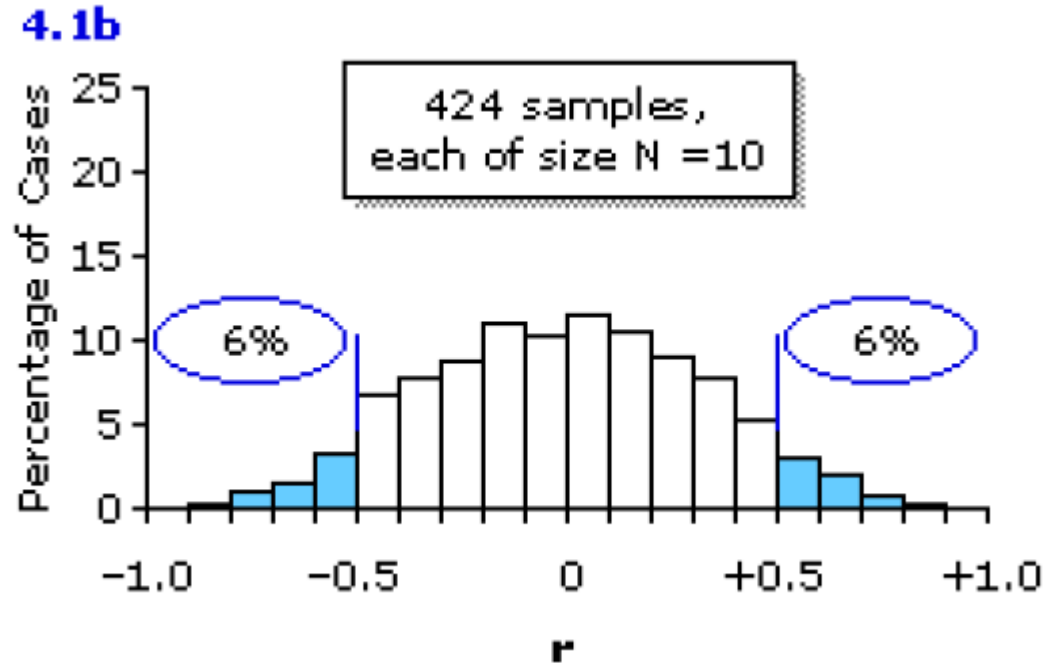
**Practice: Linear Correlation Coefficient for 5 throws of a pair of dice, done 473 times**

From Richard Lowry's class notes:

<http://www.vassarstats.net/textbook/ch4pt1.html>

# AOSC 652: Analysis Methods in AOSC

## Test for statistical significance of correlation vital



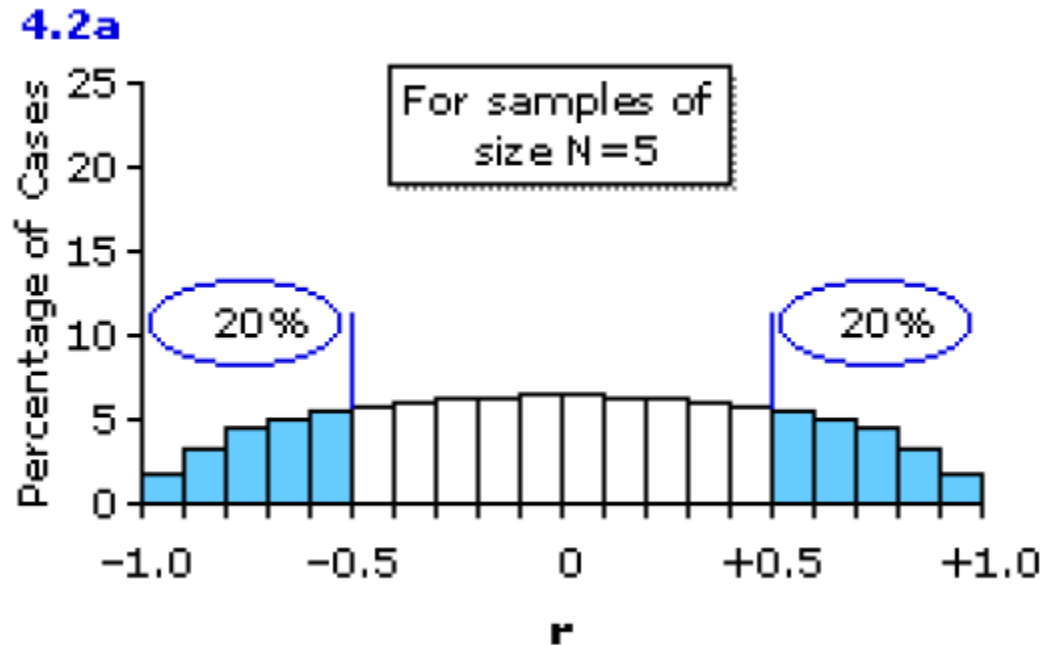
**Practice: Linear Correlation Coefficient for 10 throws of a pair of dice, done 424 times**

From Richard Lowry's class notes:

<http://www.vassarstats.net/textbook/ch4pt1.html>

# AOSC 652: Analysis Methods in AOSC

## Test for statistical significance of correlation vital



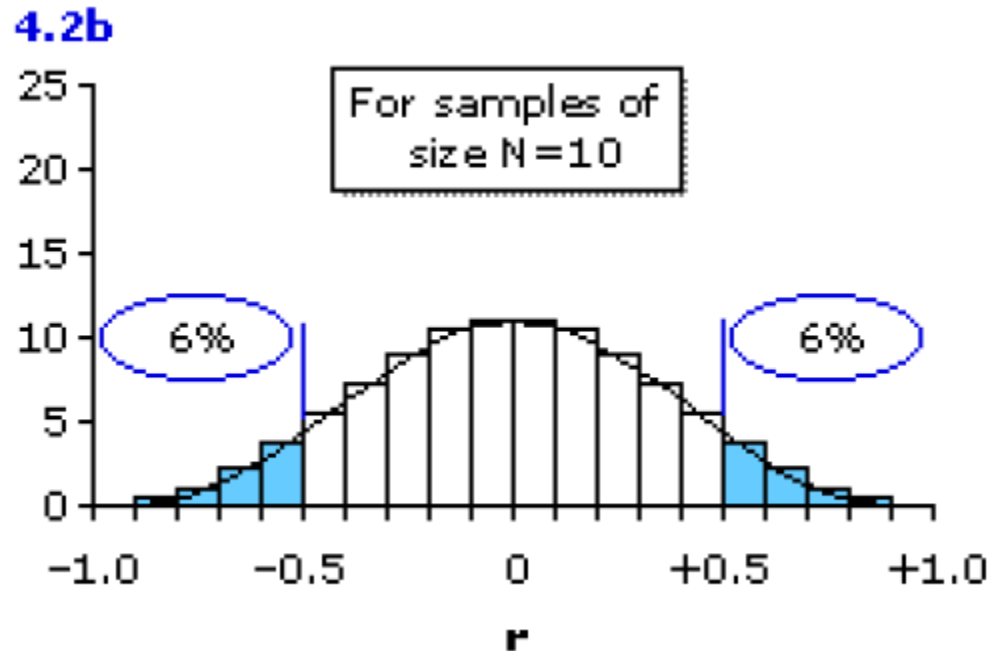
**Theory: Correlation Coefficient for 5 throws of a pair of dice, done infinite # of times**

From Richard Lowry's class notes:

<http://www.vassarstats.net/textbook/ch4pt1.html>

# AOSC 652: Analysis Methods in AOSC

## Test for statistical significance of correlation vital



**Theory: Correlation Coefficient for 10 throws of a pair of dice, done infinite # of times**

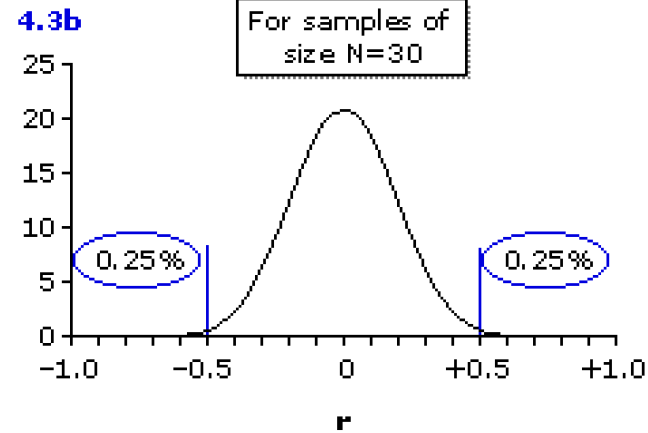
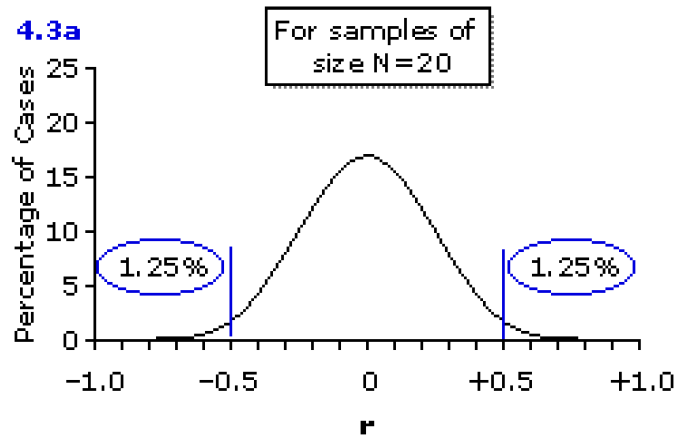
From Richard Lowry's class notes:

<http://www.vassarstats.net/textbook/ch4pt1.html>



# AOSC 652: Analysis Methods in AOSC

## Test for statistical significance of correlation vital



**Theory: Correlation Coefficient for 20 throws (left) and 30 throws (right) of a pair of dice (left)**

From Richard Lowry's class notes:

<http://www.vassarstats.net/textbook/ch4pt1.html>

# AOSC 652: Analysis Methods in AOSC

## Test for statistical significance of correlation vital

**Absolute value of  $r$  needed to show two quantities bear a statistically significant relation at the 95% confidence interval, as a function of sample size**

<b>N</b>	<b>Directional</b>	<b>Non-Directional</b>
5	0.81	0.88

From Richard Lowry's class notes:

<http://www.vassarstats.net/textbook/ch4pt1.html>

# AOSC 652: Analysis Methods in AOSC

## Test for statistical significance of correlation vital

**Absolute value of  $r$  needed to show two quantities bear a statistically significant relation at the 95% confidence interval, as a function of sample size**

<b>N</b>	<b>Directional</b>	<b>Non-Directional</b>
5	0.81	0.88
10	0.55	0.63
15	0.44	0.51
20	0.38	0.44
30	0.30	0.35

From Richard Lowry's class notes:

<http://www.vassarstats.net/textbook/ch4pt1.html>

# AOSC 652: Analysis Methods in AOSC

## Test for statistical significance of correlation vital

**Absolute value of  $r$  needed to show two quantities bear a statistically significant relation at the 95% confidence interval, as a function of sample size**

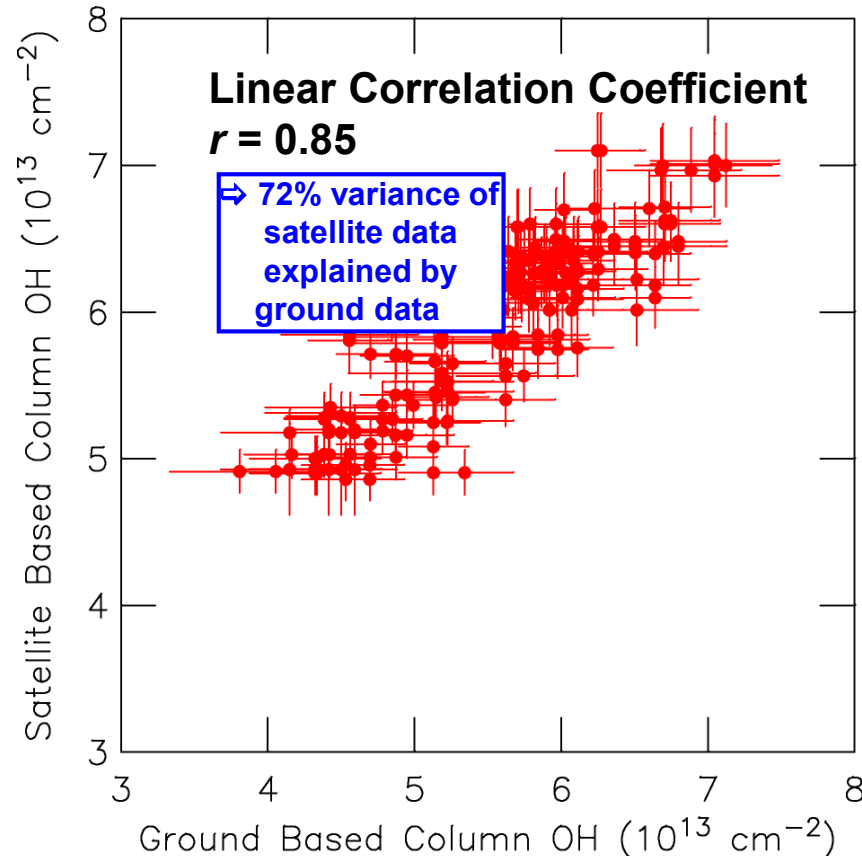
<b>N</b>	<b>Directional</b>	<b>Non-Directional</b>
5	0.81	0.88
10	0.55	0.63
15	0.44	0.51
20	0.38	0.44
30	0.30	0.35

**See Section 8.4 of von Storch and Zwiers, *Statistical Analysis of Climate Research* for methodology for assessing statistical significance of a linear correlation**

**<http://www.leif.org/EOS/vonSt0521012309.pdf>**

# AOSC 652: Analysis Methods in AOSC

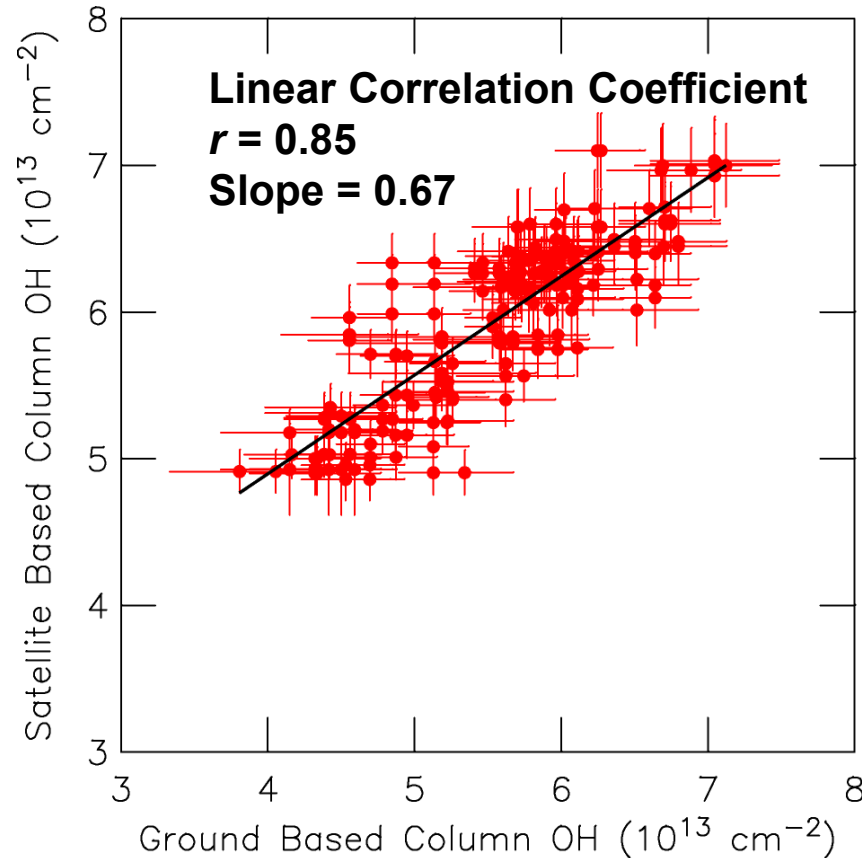
Suppose you have two sets of measurements (or data and model) that you'd like to relate.



What else might we want to know?

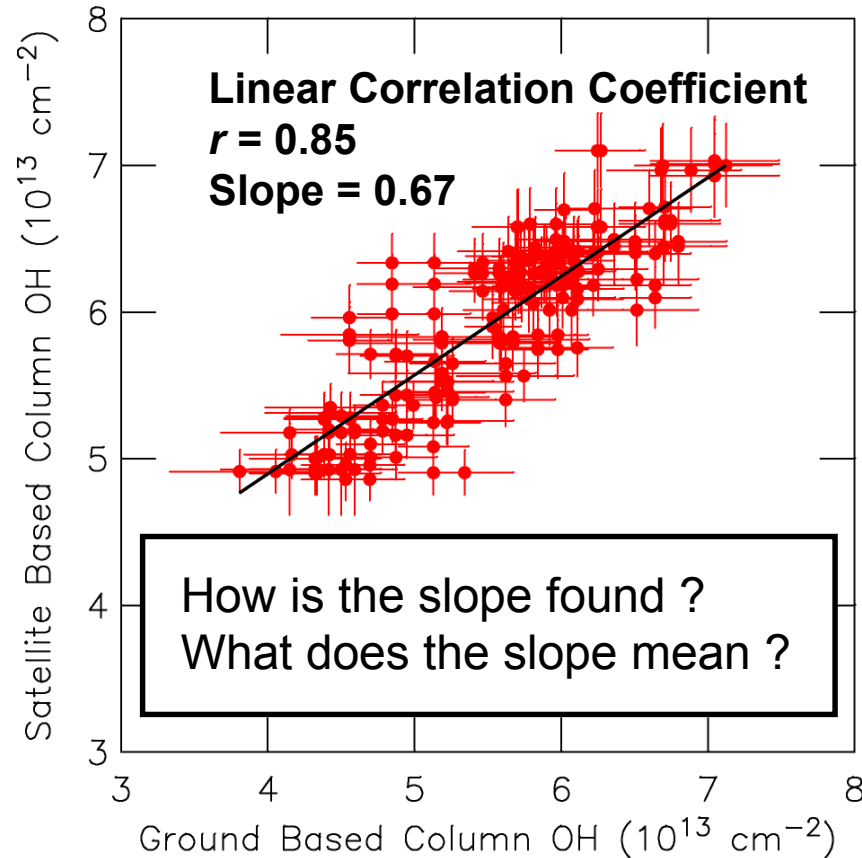
# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



# AOSC 652: Analysis Methods in AOSC

## Linear Least Squares Fitting:

$$y = a + b x$$

## Minimize:

$$\sum_{i=1}^N (a + b x_i - y_i)^2 \equiv \text{Cost Function}$$

$$\frac{\partial \text{Cost Function}}{\partial a} = 2 \sum_{i=1}^N (a + b x_i - y_i) = 0$$

$$\frac{\partial \text{Cost Function}}{\partial b} = 2 \sum_{i=1}^N (a + b x_i - y_i) x_i = 0$$



# AOSC 652: Analysis Methods in AOSC

$$\frac{\partial \text{Cost Function}}{\partial a} = 2 \sum_{i=1}^N (a + b x_i - y_i) = 0$$

$$\frac{\partial \text{Cost Function}}{\partial b} = 2 \sum_{i=1}^N (a + b x_i - y_i) x_i = 0$$

$$a N + b \sum x_i - \sum y_i = 0$$

$$a \sum x_i + b \sum x_i^2 - \sum x_i \cdot y_i = 0$$

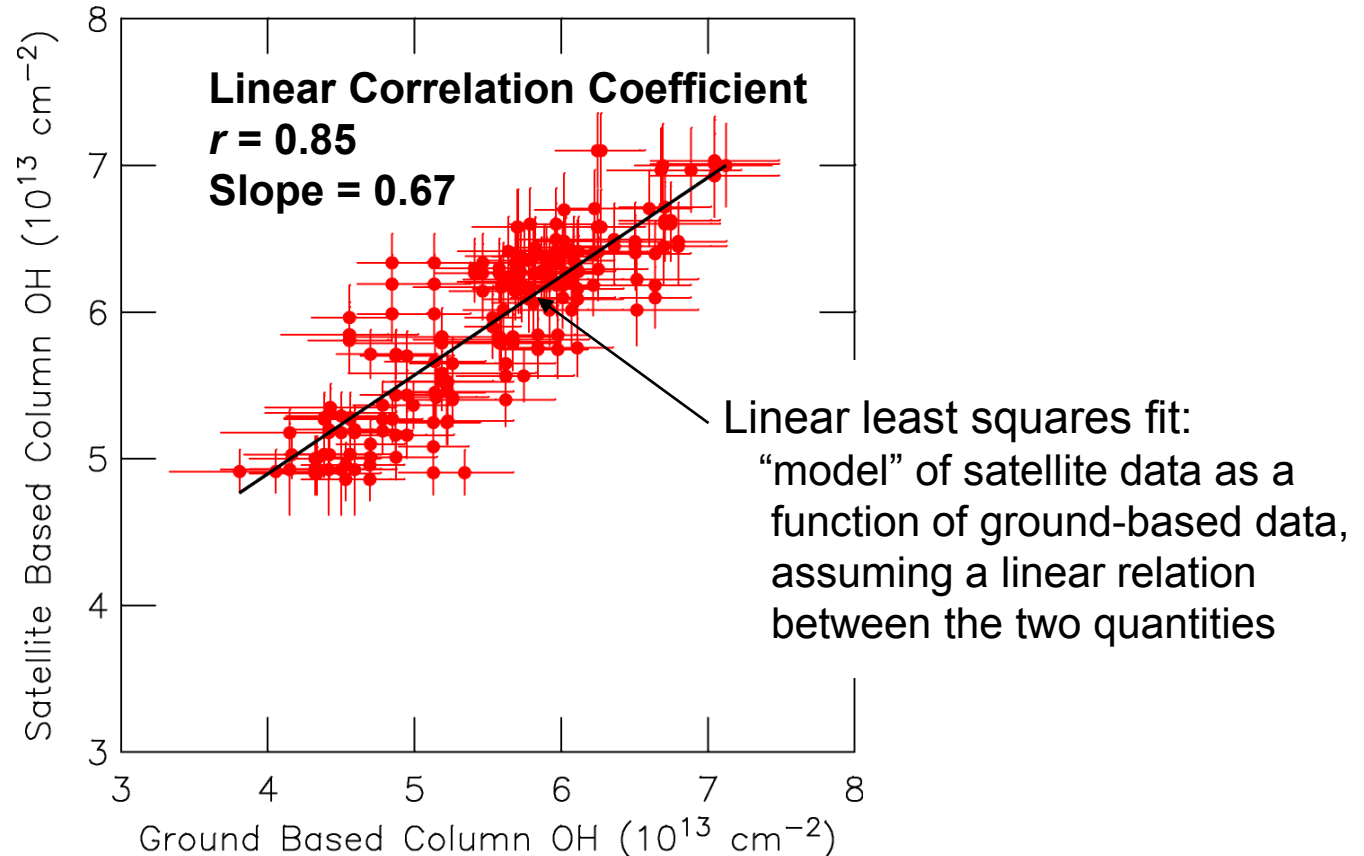
Can show:

$$b = \frac{\sum x_i \cdot \sum y_i - N \cdot \sum x_i y_i}{(\sum x_i)^2 - N \cdot \sum x_i^2} \quad \Leftrightarrow \text{Slope of linear least squares fit}$$

$$a = \frac{\sum y_i - b \cdot \sum x_i}{N} \quad \Leftrightarrow \text{Intercept of linear least squares fit}$$

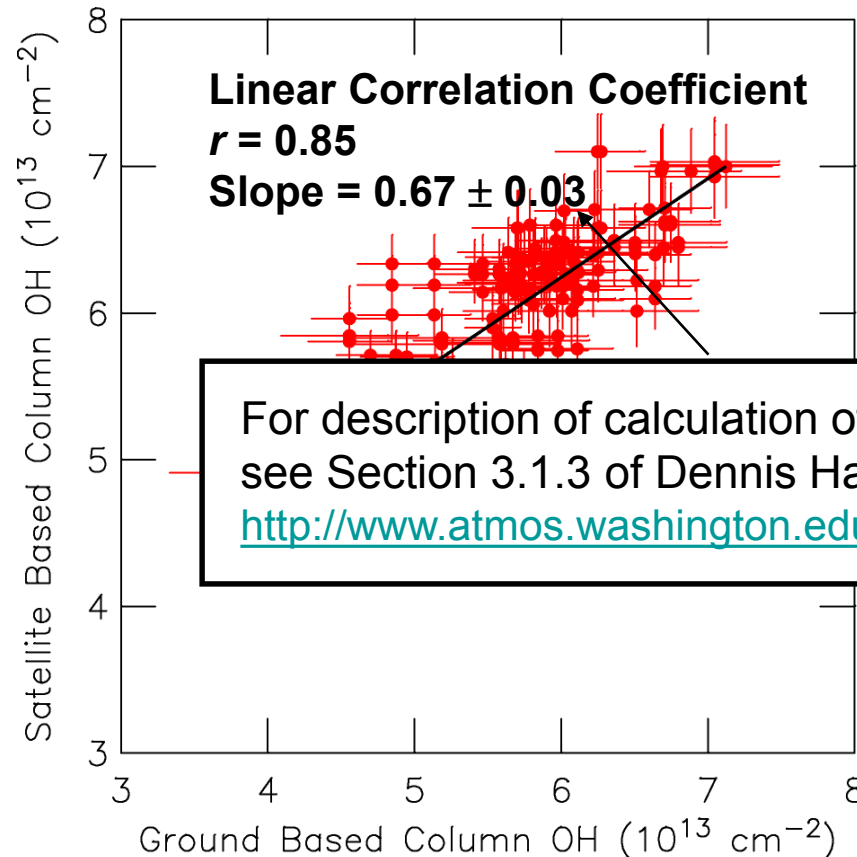
# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



# AOSC 652: Analysis Methods in AOSC

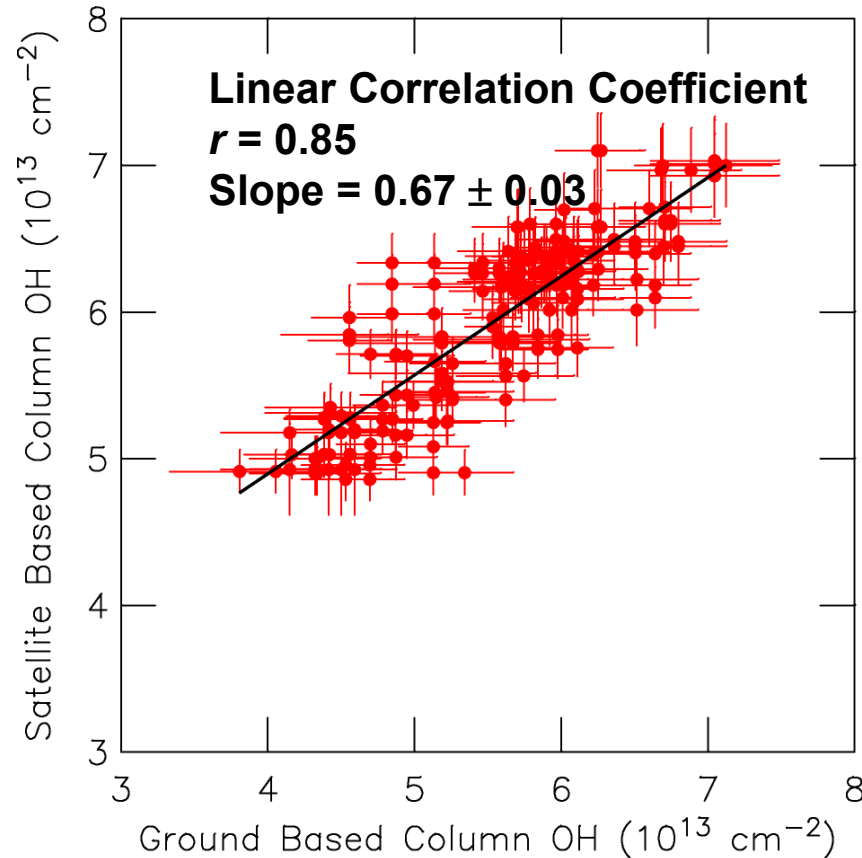
Suppose you have two sets of measurements (or data and model) that you'd like to relate.



For description of calculation of the uncertainty of slope see Section 3.1.3 of Dennis Hartmann's class notes:  
[http://www.atmos.washington.edu/~dennis/552\\_Notes\\_3.pdf](http://www.atmos.washington.edu/~dennis/552_Notes_3.pdf)

# AOSC 652: Analysis Methods in AOSC

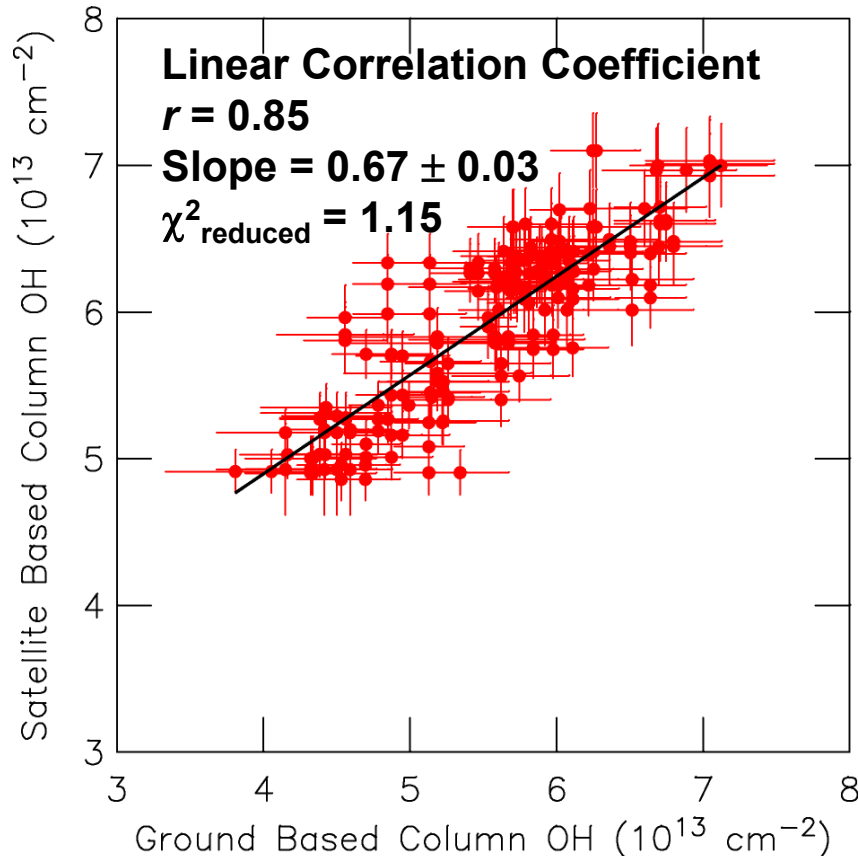
Suppose you have two sets of measurements (or data and model) that you'd like to relate.



How can we quantify the “goodness of the fit” ?

# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



Reduced  $\chi^2$  :

$$\chi^2_{\text{Fit}} = \frac{1}{\nu} \sum_{i=1}^N \left( \frac{y_i - \text{Fit}(x_i)}{\sigma_i} \right)^2$$

where  $\sigma_i$  is the uncertainty associated with each measurement and

$$\nu = N - p - 1$$

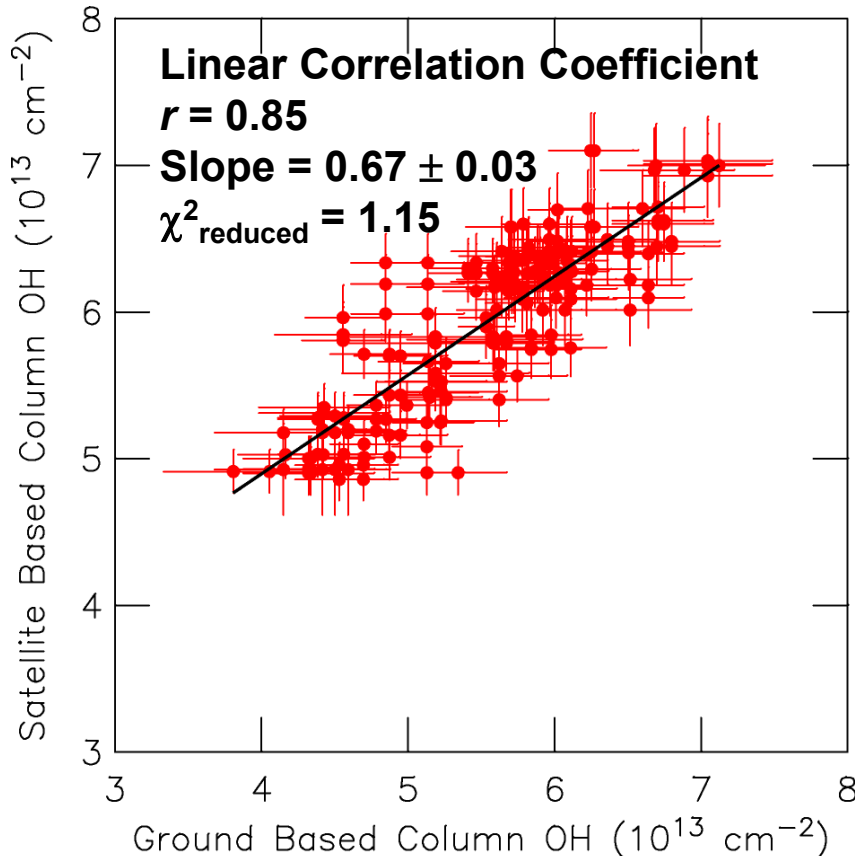
for  $N$  = number of observations

$p$  = number of fitting parameters

Often, one sets  $\nu$  equal to  $N$

# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



Reduced  $\chi^2$  :

$$\chi^2_{\text{Fit}} = \frac{1}{\nu} \sum_{i=1}^N \left( \frac{y_i - \text{Fit}(x_i)}{\sigma_i} \right)^2$$

where  $\sigma_i$  is the uncertainty associated with each measurement and

$$\nu = N - p - 1$$

for  $N$  = number of observations

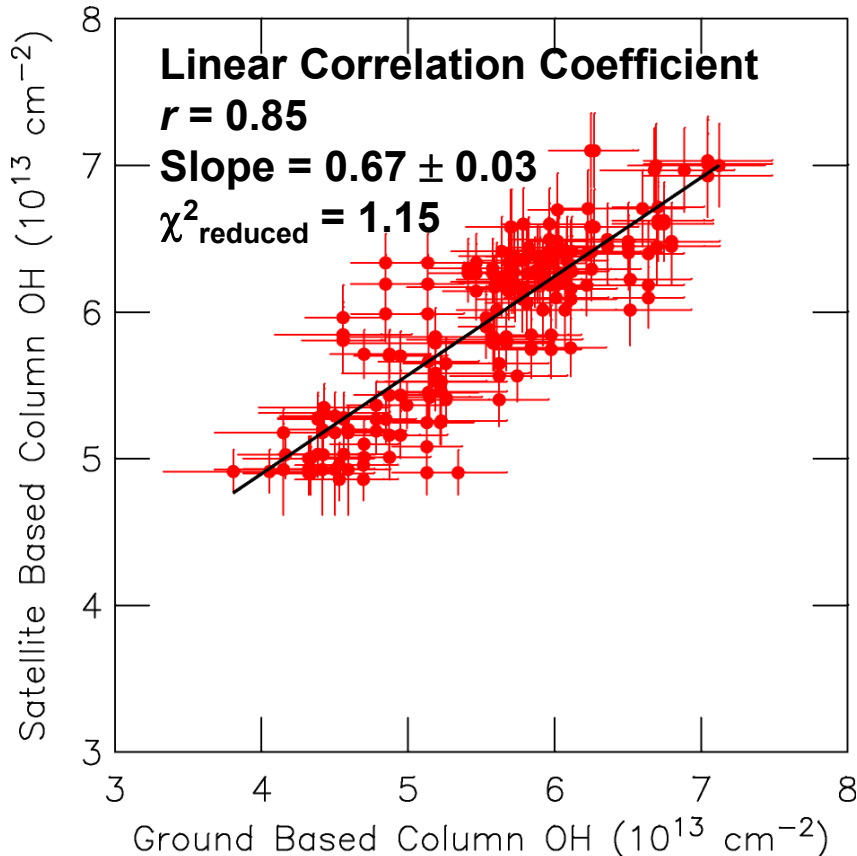
$p$  = number of fitting parameters

Often, one sets  $\nu$  equal to  $N$

Degrees of freedom

# AOSC 652: Analysis Methods in AOSC

Suppose you have two sets of measurements (or data and model) that you'd like to relate.



Reduced  $\chi^2$  :

$$\chi^2_{\text{Fit}} = \frac{1}{\nu} \sum_{i=1}^N \left( \frac{y_i - \text{Fit}(x_i)}{\sigma_i} \right)^2$$

where  $\sigma_i$  is the uncertainty associated with each measurement and

$$\nu = N - p - 1$$

for  $N$  = number of observations

$p$  = number of fitting parameters

Often, one sets  $\nu$  equal to  $N$

**Reduced  $\chi^2$  is an under-used, often unappreciated, quantitative measure of the degree of relation between a model and data, or two types of measurements (or, two types of models, if we can define uncertainties for each)**

# AOSC 652: Analysis Methods in AOSC

**Reduced  $\chi^2$  is an under-used, perhaps unappreciated quantitative measure of the degree of relation between a model and data, or two types of measurements (or, two types of models, if we can define uncertainties for each)**



# AOSC 652: Analysis Methods in AOSC

Reduced  $\chi^2$  is an under-used, perhaps unappreciated quantitative measure of the degree of relation between a model and data, or two types of measurements (or, two types of models, if we can define uncertainties for each)

**Reduced  $\chi^2$  is very commonly used in the physics community**

**Indeed, in physics lab courses, students are sometimes cautioned to be critical of an experiment where reduced  $\chi^2$  lies close to zero.**

**Why is this?**

# AOSC 652: Analysis Methods in AOSC

Reduced  $\chi^2$  is an under-used, perhaps unappreciated quantitative measure of the degree of relation between a model and data, or two types of measurements (or, two types of models, if we can define uncertainties for each)

Reduced  $\chi^2$  is very commonly used in the physics community

Indeed, in physics lab courses, students are sometimes cautioned to be critical of an experiment where reduced  $\chi^2$  lies close to zero.

Why is this?

See [http://www.physics.csbsju.edu/stats/chi\\_fit.html](http://www.physics.csbsju.edu/stats/chi_fit.html) for more for reduced  $\chi^2$

# A method for evaluating bias in global measurements of CO<sub>2</sub> total columns from space

D. Wunch<sup>1</sup>, P. O. Wennberg<sup>1</sup>, G. C. Toon<sup>1,2</sup>, B. J. Connor<sup>3</sup>, B. Fisher<sup>2</sup>, G. B. Osterman<sup>2</sup>, C. Frankenberg<sup>2</sup>, L. Mandrake<sup>2</sup>, C. O'Dell<sup>4</sup>, P. Ahonen<sup>5</sup>, S. C. Biraud<sup>14</sup>, R. Castano<sup>2</sup>, N. Cressie<sup>6</sup>, D. Crisp<sup>2</sup>, N. M. Deutscher<sup>7,8</sup>, A. Eldering<sup>2</sup>, M. L. Fisher<sup>14</sup>, D. W. T. Griffith<sup>8</sup>, M. Gunson<sup>2</sup>, P. Heikkinen<sup>5</sup>, G. Keppel-Aleks<sup>1</sup>, E. Kyrö<sup>5</sup>, R. Lindenmaier<sup>15</sup>, R. Macatangay<sup>8</sup>, J. Mendonca<sup>15</sup>, J. Messerschmidt<sup>7</sup>, C. E. Miller<sup>2</sup>, I. Morino<sup>9</sup>, J. Notholt<sup>7</sup>, F. A. Oyafuso<sup>2</sup>, M. Rettinger<sup>10</sup>, J. Robinson<sup>12</sup>, C. M. Roehl<sup>1</sup>, R. J. Salawitch<sup>11</sup>, V. Sherlock<sup>12</sup>, K. Strong<sup>15</sup>, R. Sussmann<sup>10</sup>, T. Tanaka<sup>9,\*</sup>, D. R. Thompson<sup>2</sup>, O. Uchino<sup>9</sup>, T. Warneke<sup>7</sup>, and S. C. Wofsy<sup>13</sup>

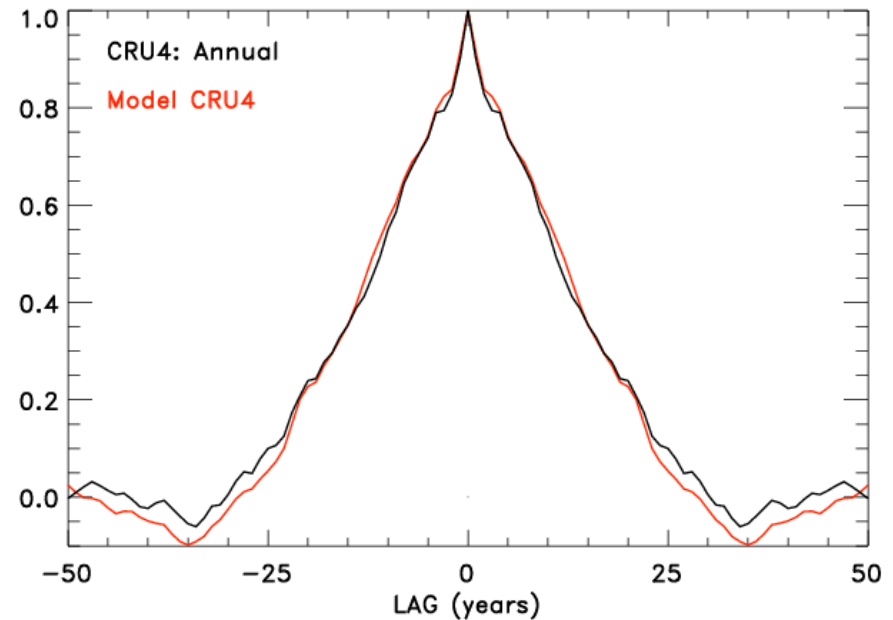
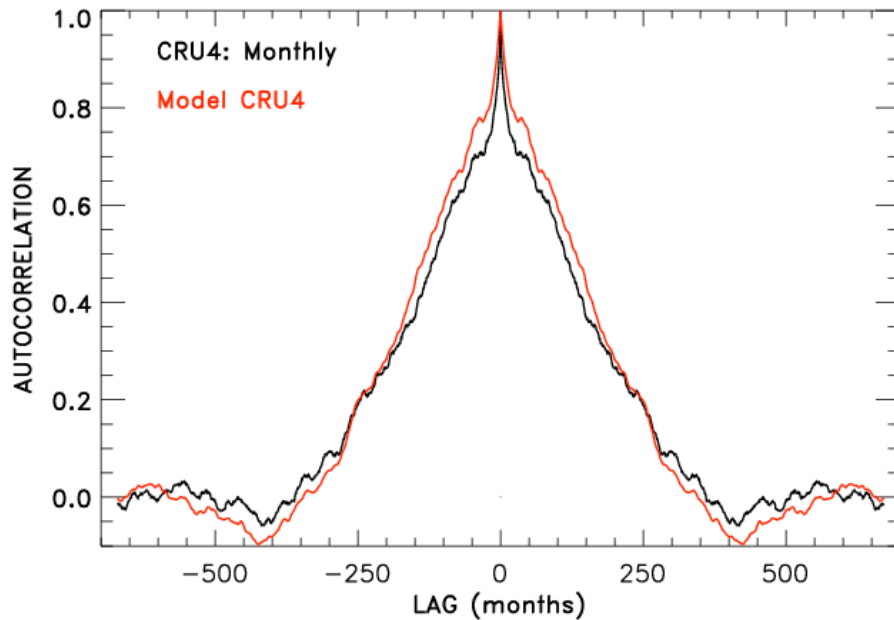
Retrievals are defined as successful by the master quality flag when they satisfy  $\chi^2 < 1.2$ . However, the  $\chi^2$  values have increased linearly over time, because the time-dependent radiometric calibration caused by a sensitivity degradation of the O<sub>2</sub> A-band channel was not applied to the noise model. To compensate for this, we adjust the cutoff value so that it starts at 1.2 and evolves with a linear increase in time, matching the increase in minimum  $\chi^2$ . As a result, a similar number of scenes are retained over time.

# AOSC 652: Analysis Methods in AOSC

## Auto-correlation

$$\phi(L) = \frac{1}{N-2L} \sum_{k=L}^{N-L} x'_k x'_{k+L} = \overline{x'_k x'_{k+L}}; L = 0, \pm 1, \pm 2, \pm 3, \dots$$

where  $x'$  is deviation from mean



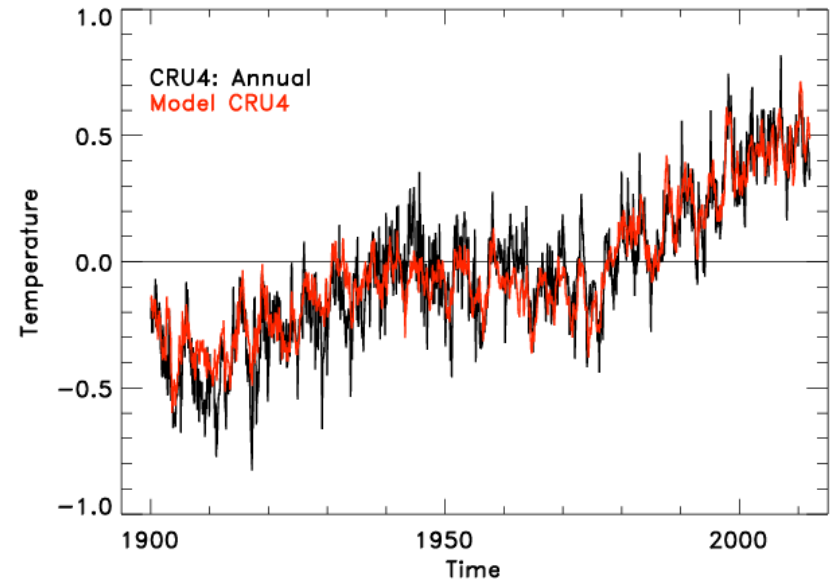
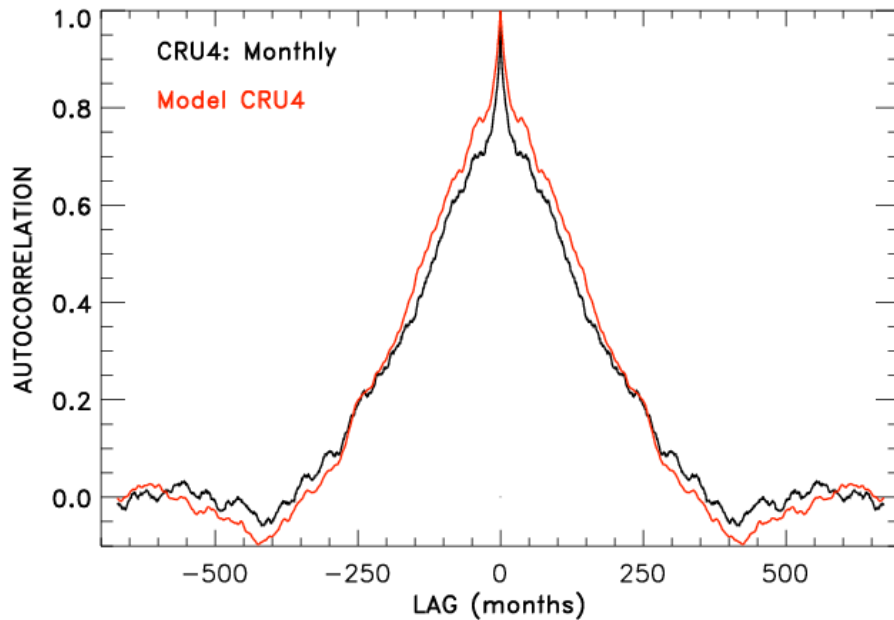
Canty et al., ACP, 2013

# AOSC 652: Analysis Methods in AOSC

## Auto-correlation

$$\phi(L) = \frac{1}{N-2L} \sum_{k=L}^{N-L} x'_k x'_{k+L} = \overline{x'_k x'_{k+L}}; L = 0, \pm 1, \pm 2, \pm 3, \dots$$

where  $x'$  is deviation from mean



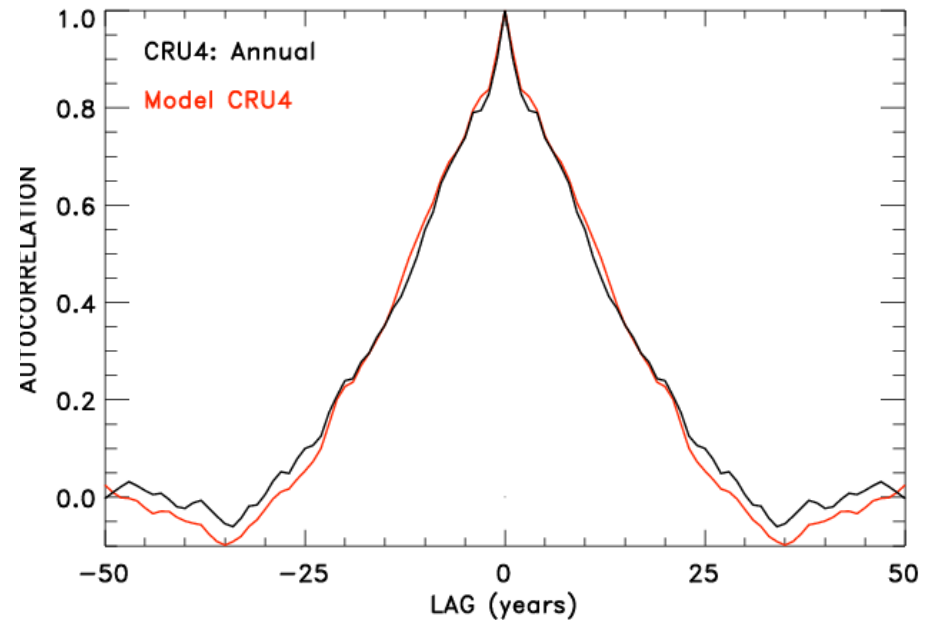
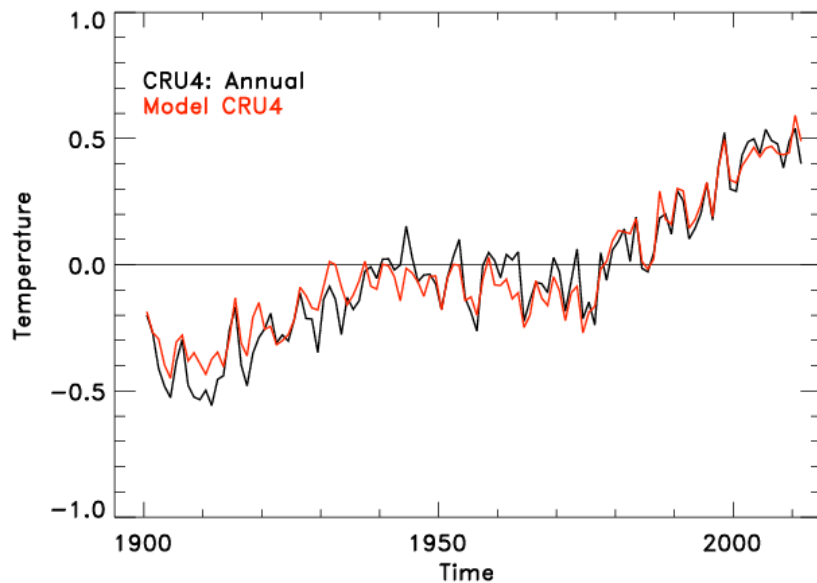
Canty et al., ACP, 2013

# AOSC 652: Analysis Methods in AOSC

## Auto-correlation

$$\phi(L) = \frac{1}{N-2L} \sum_{k=L}^{N-L} x'_k x'_{k+L} = \overline{x'_k x'_{k+L}}; L = 0, \pm 1, \pm 2, \pm 3, \dots$$

where  $x'$  is deviation from mean



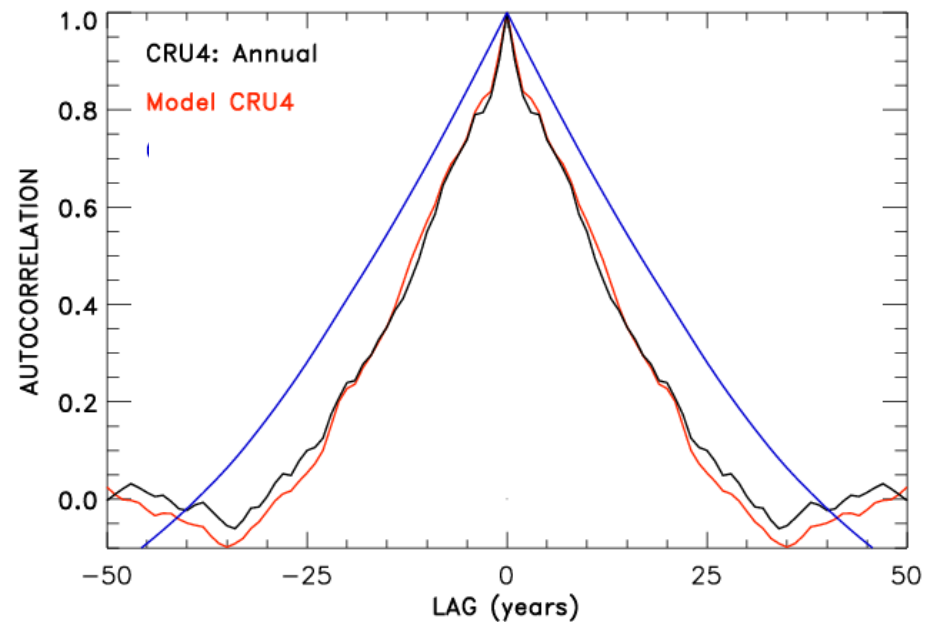
Canty et al., ACP, 2013

# AOSC 652: Analysis Methods in AOSC

Auto-correlation

$$\phi(L) = \frac{1}{N-2L} \sum_{k=L}^{N-L} x'_k x'_{k+L} = \overline{x'_k x'_{k+L}}; L = 0, \pm 1, \pm 2, \pm 3, \dots$$

where  $x'$  is deviation from mean



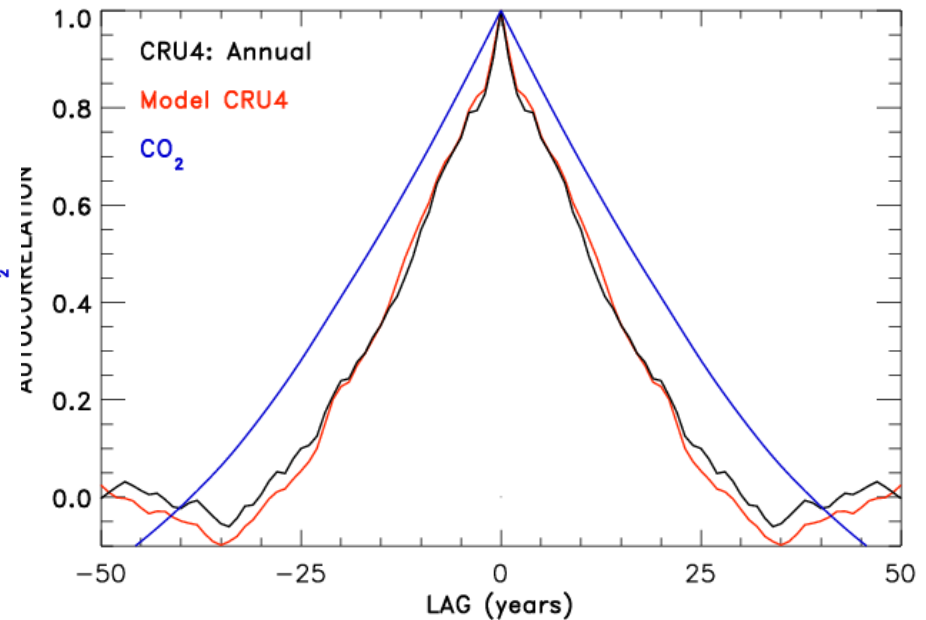
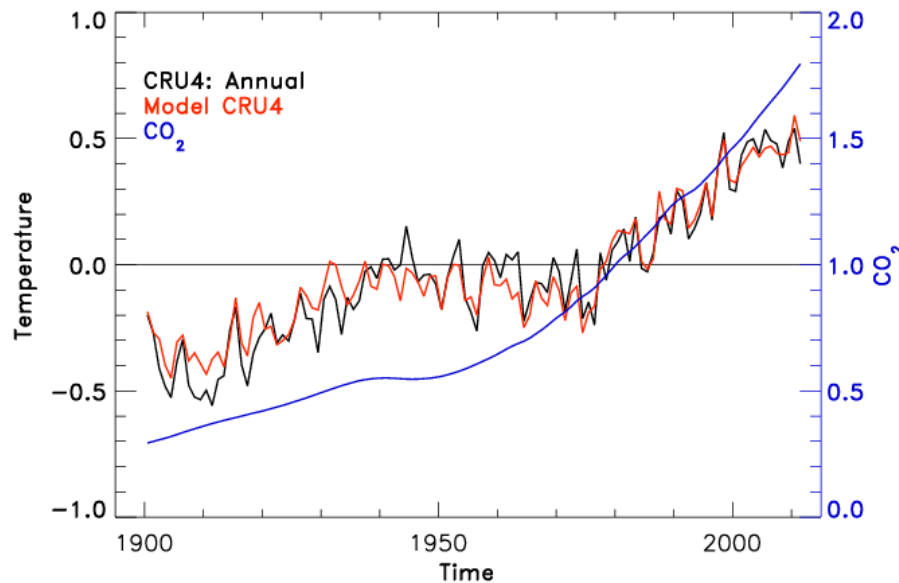
Canty et al., ACP, 2013

# AOSC 652: Analysis Methods in AOSC

## Auto-correlation

$$\phi(L) = \frac{1}{N-2L} \sum_{k=L}^{N-L} x'_k x'_{k+L} = \overline{x'_k x'_{k+L}}; \quad L = 0, \pm 1, \pm 2, \pm 3, \dots$$

where  $x'$  is deviation from mean



Canty et al., ACP, 2013

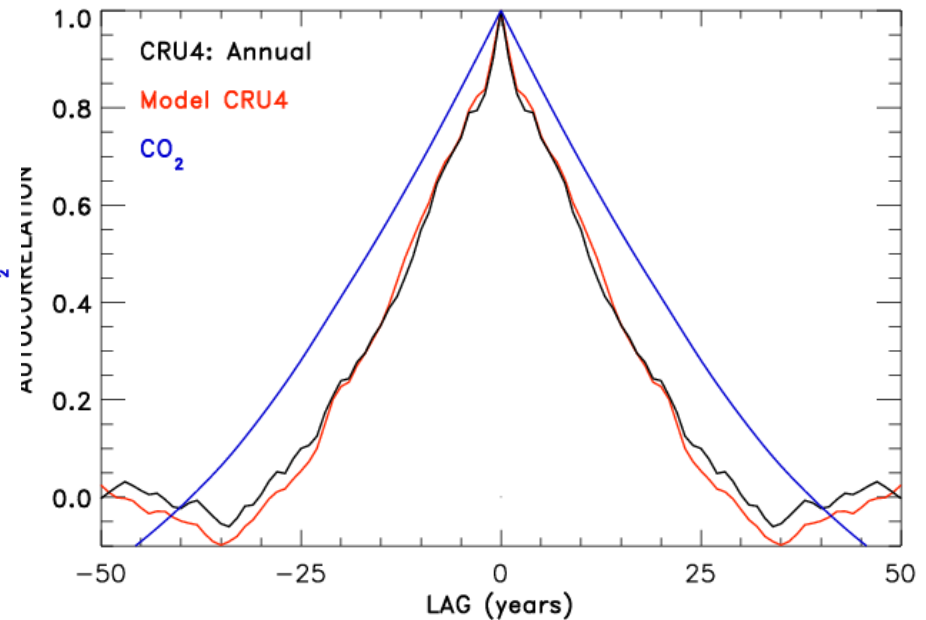
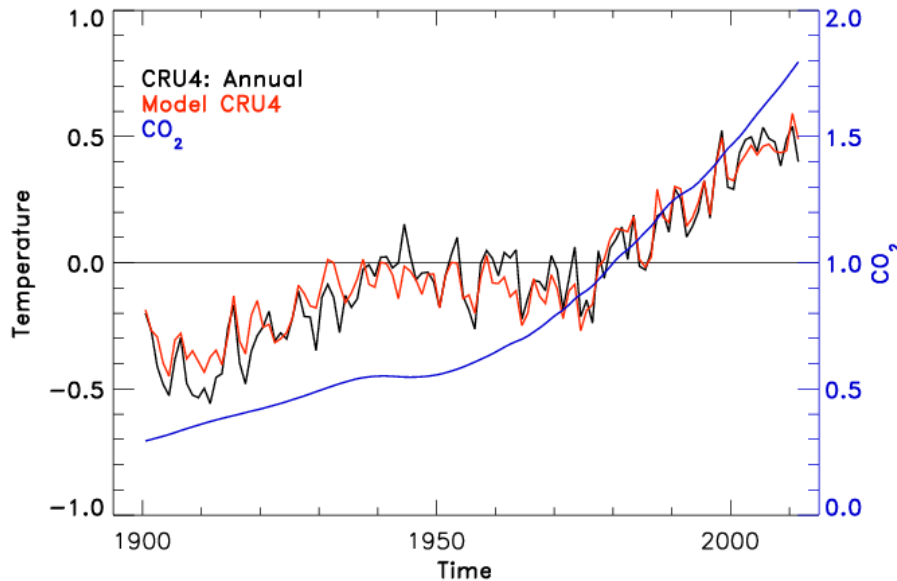


# AOSC 652: Analysis Methods in AOSC

## Auto-correlation

$$\phi(L) = \frac{1}{N-2L} \sum_{k=L}^{N-L} x'_k x'_{k+L} = \overline{x'_k x'_{k+L}}; L = 0, \pm 1, \pm 2, \pm 3, \dots$$

where  $x'$  is deviation from mean



Canty et al., ACP, 2013

**The behavior of the auto-correlation of a signal is often used to infer degrees of freedom: see for example Dennis Hartmann's class notes:**  
[http://www.atmos.washington.edu/~dennis/552\\_Notes\\_6a.pdf](http://www.atmos.washington.edu/~dennis/552_Notes_6a.pdf)