

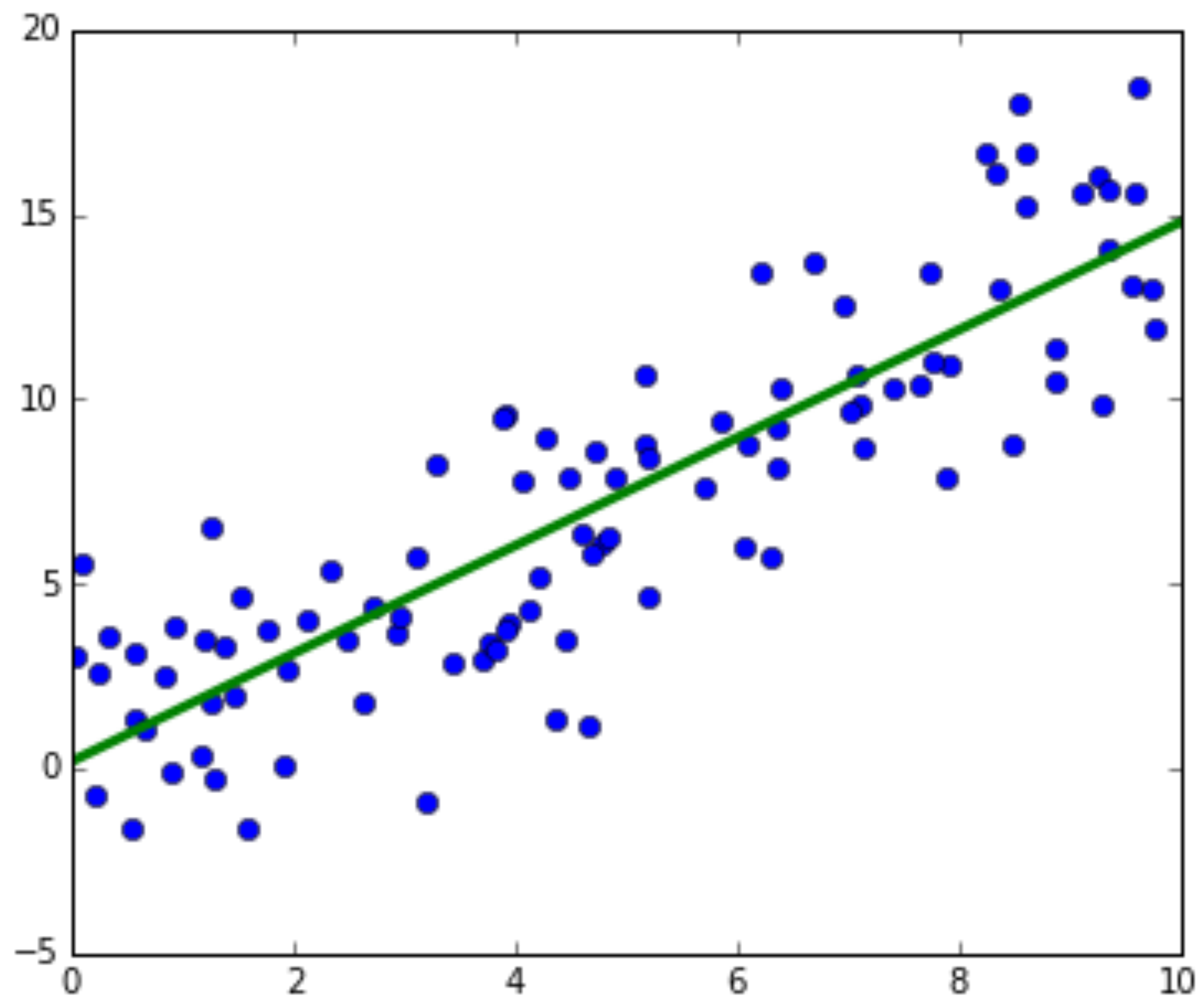
Info 2950, Lecture 22

25 Apr 2017

Prob Set 6: due Mon night 24 Apr

Prob Set 7: to be issued tonight, due mid next week

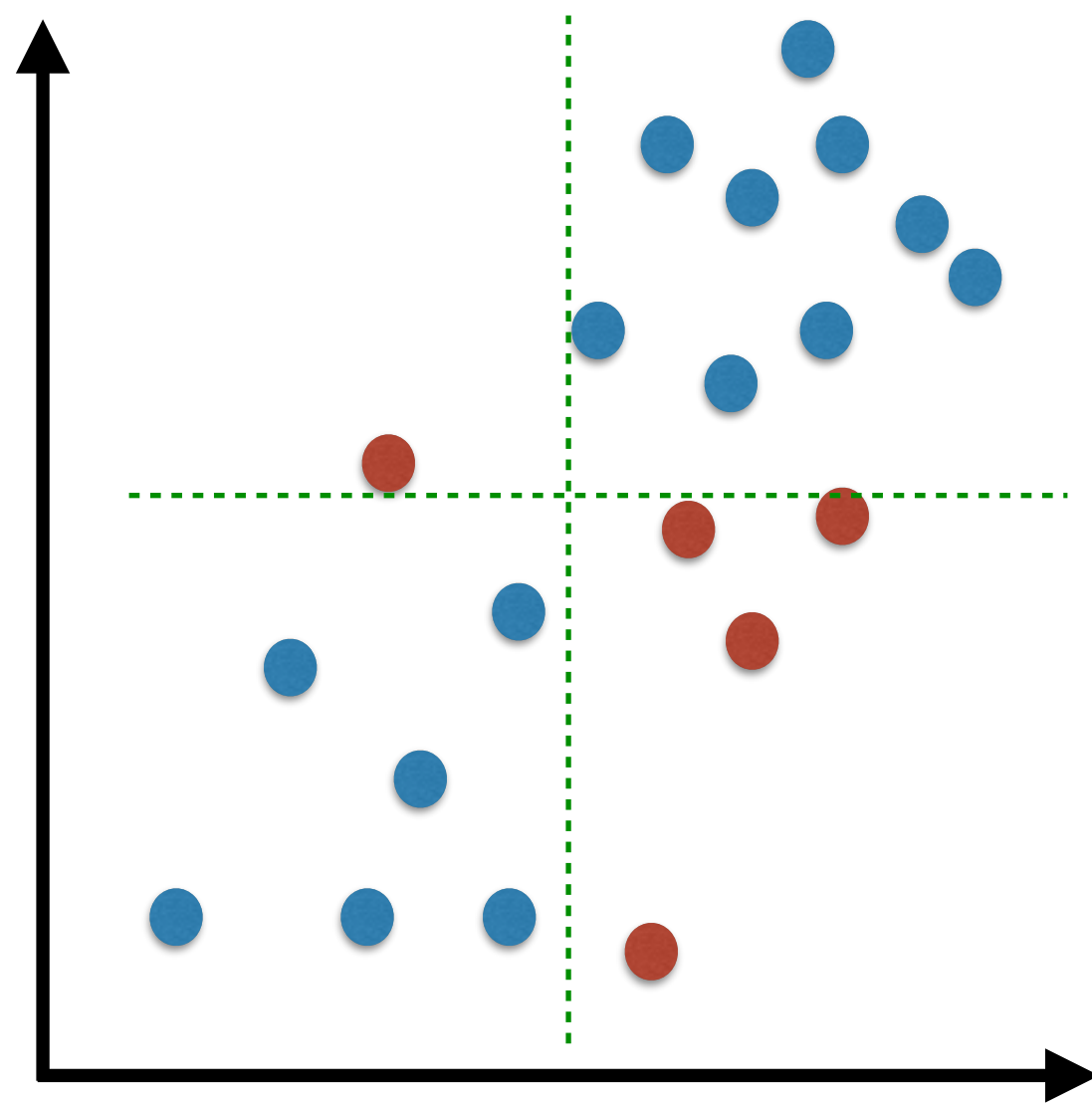
Prob Set 8: due 11 May (end of classes)



$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

$$\text{Cov}(X, X) = \text{Var}(X)$$

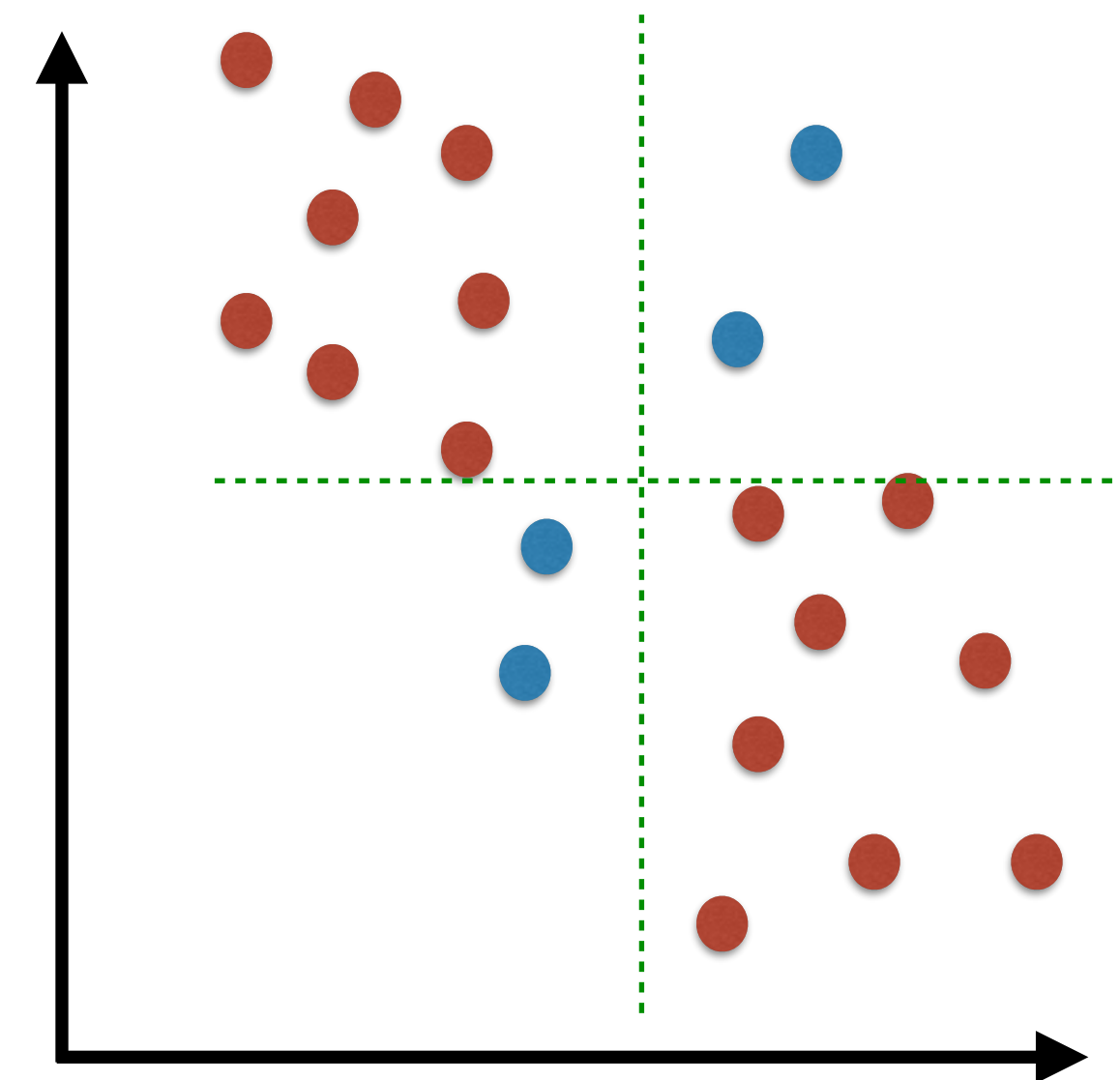
$X > E[X], Y > E[Y]$



$X < E[X], Y < E[Y]$

$\text{Cov}(X, Y) > 0$

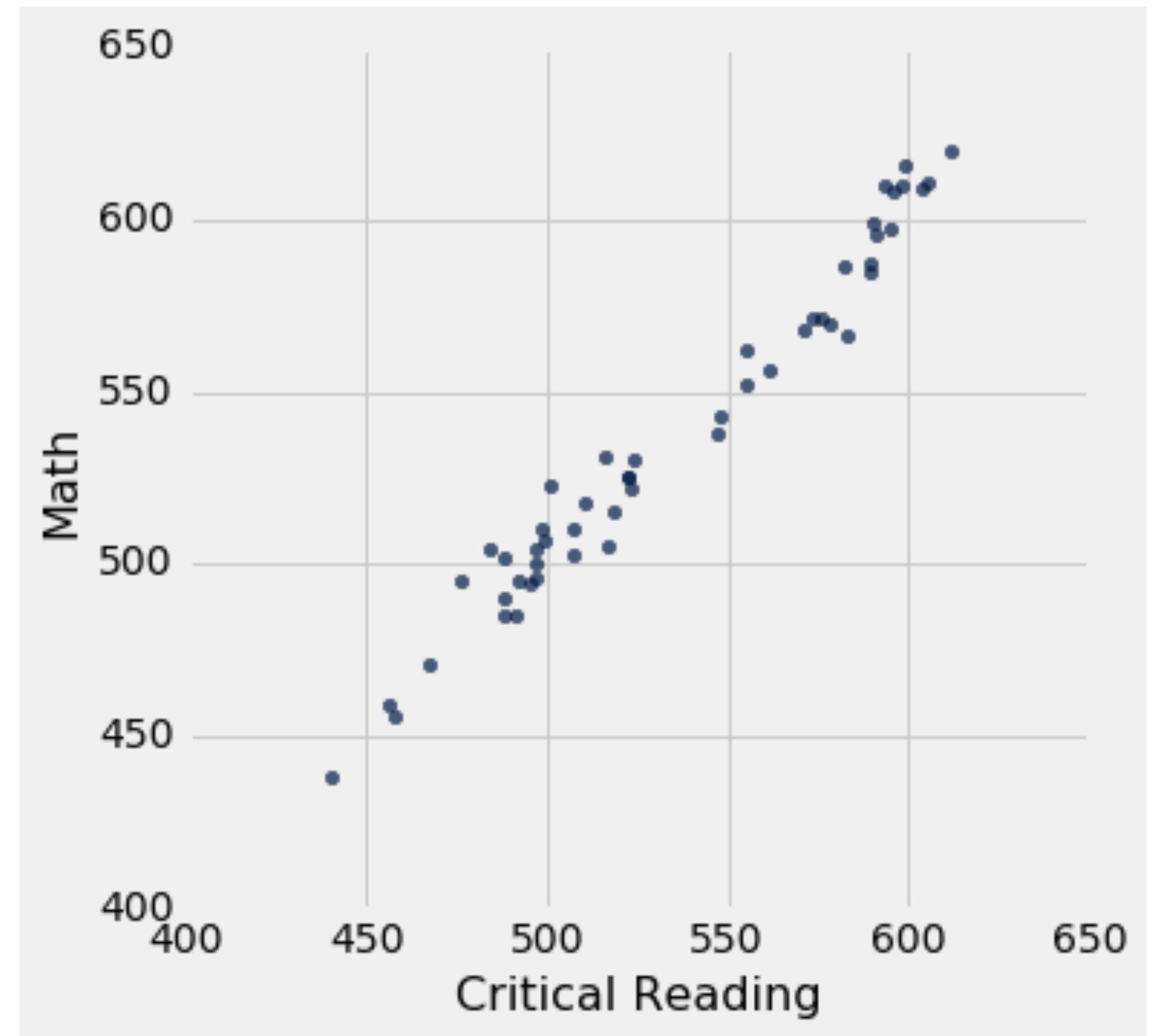
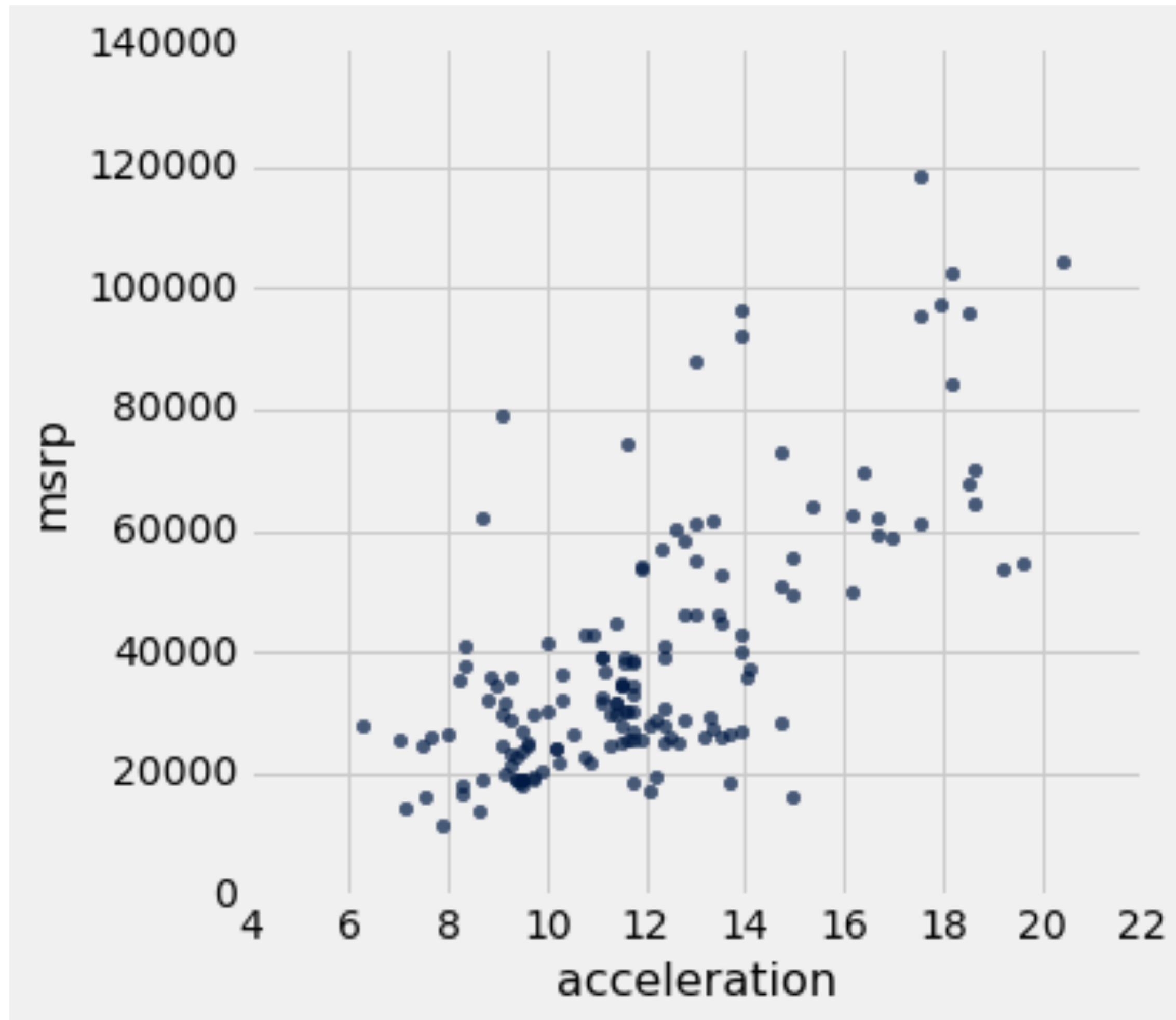
$X < E[X], Y > E[Y]$



$X > E[X], Y < E[Y]$

$\text{Cov}(X, Y) < 0$

<https://www.inferentialthinking.com/chapters/13/1/correlation.html>



<https://www.inferentialthinking.com/chapters/13/1/correlation.html>

Correlation coefficient only measures association.

Correlation does not imply causation.

Though the correlation between the weight and the math ability of children in a school district may be positive, that does not mean that doing math makes children heavier or that putting on weight improves the children's math skills.

Age is a confounding variable: older children are both heavier and better at math than younger children, on average

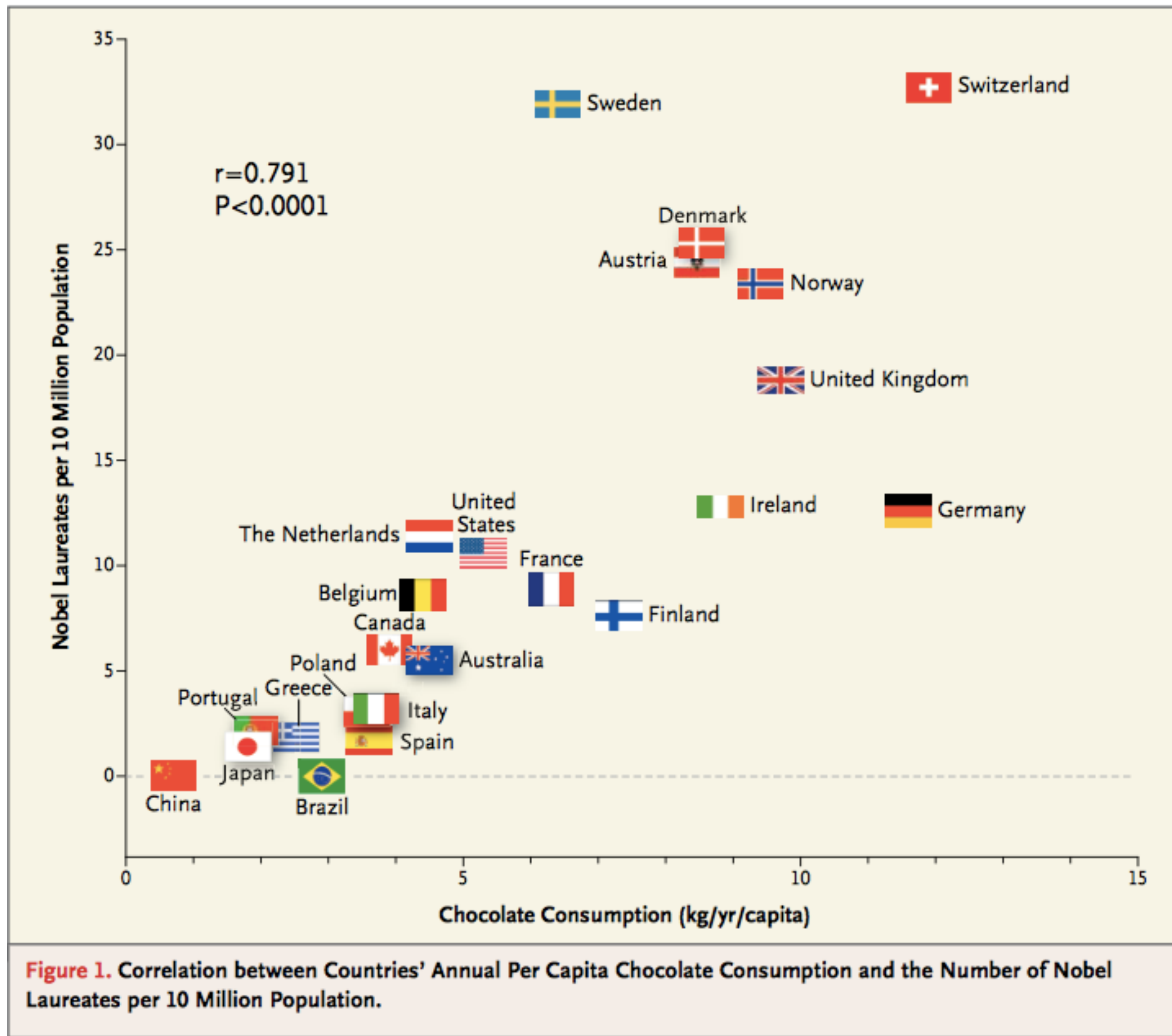


Figure 1. Correlation between Countries' Annual Per Capita Chocolate Consumption and the Number of Nobel Laureates per 10 Million Population.

see also <http://www.tylervigen.com/spurious-correlations>

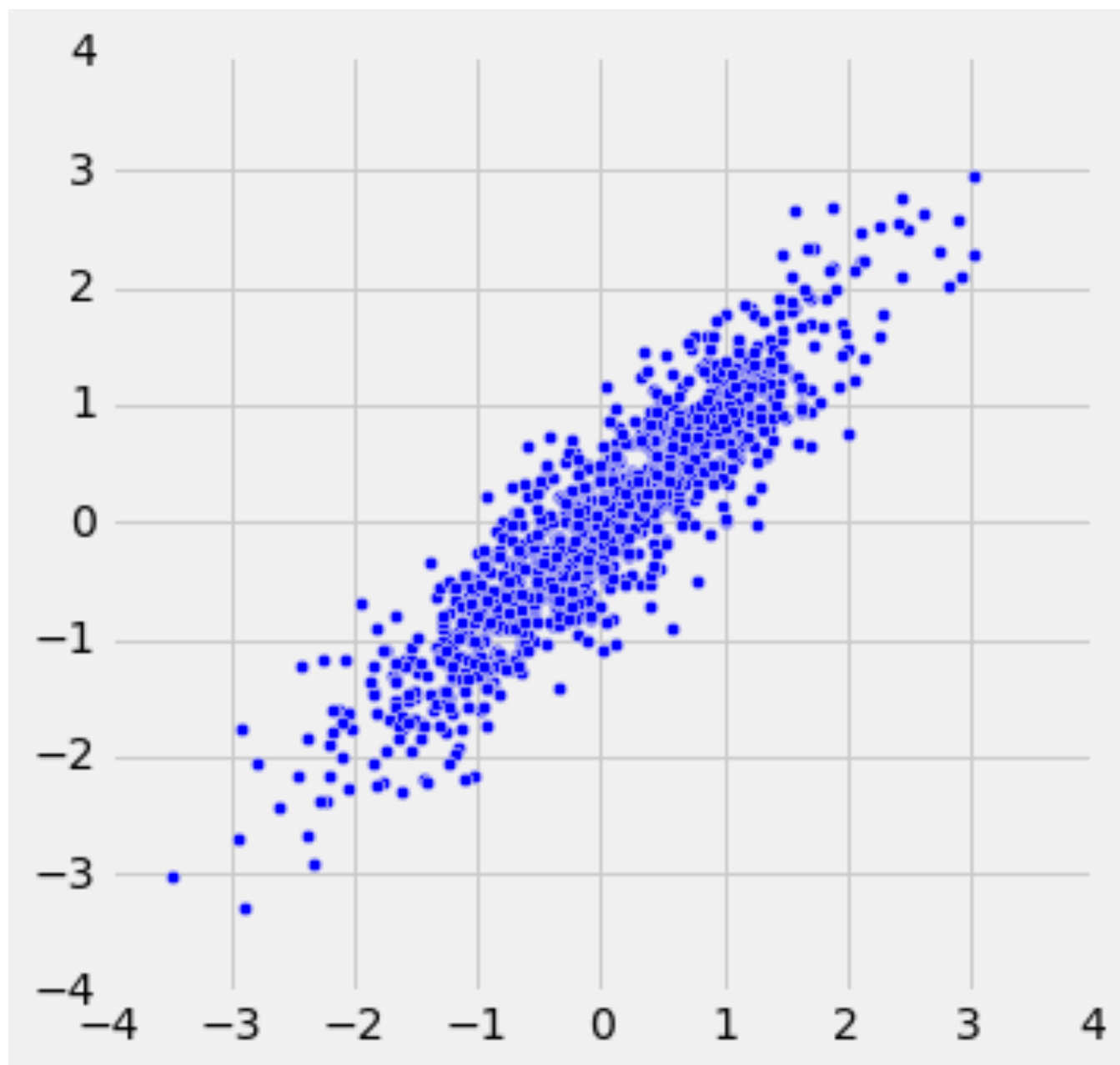
Recall that if we rescale the data, $x \rightarrow x/c$, then that divides the standard deviation $\sigma[x]$ by the same c .

In particular, we can divide by $c = \sigma[x]$, which gives a new distribution with standard deviation normalized to 1 (as is done to calculate z values).

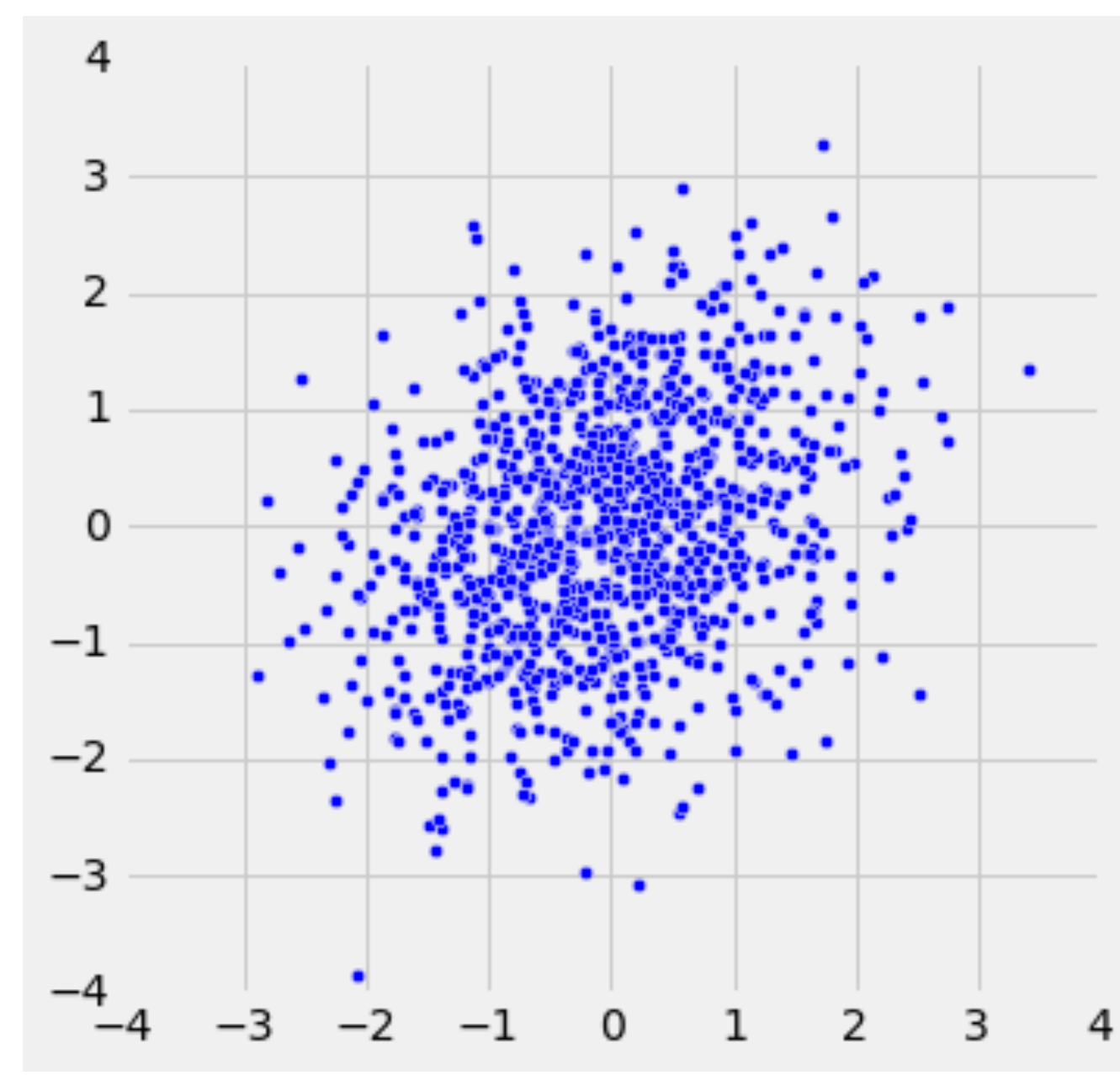
Note that $\text{Pearsonr}(x,y) = \text{Cov}[x,y] / \sigma[x] \sigma[y]$ is unchanged by rescaling both x and y .

On the other hand, if we rescale $x \rightarrow x/\sigma[x]$ and $y \rightarrow y/\sigma[y]$, then the linear regression slope $\text{Cov}[x,y]/\text{Var}[x] \rightarrow \text{Cov}[x,y]$.

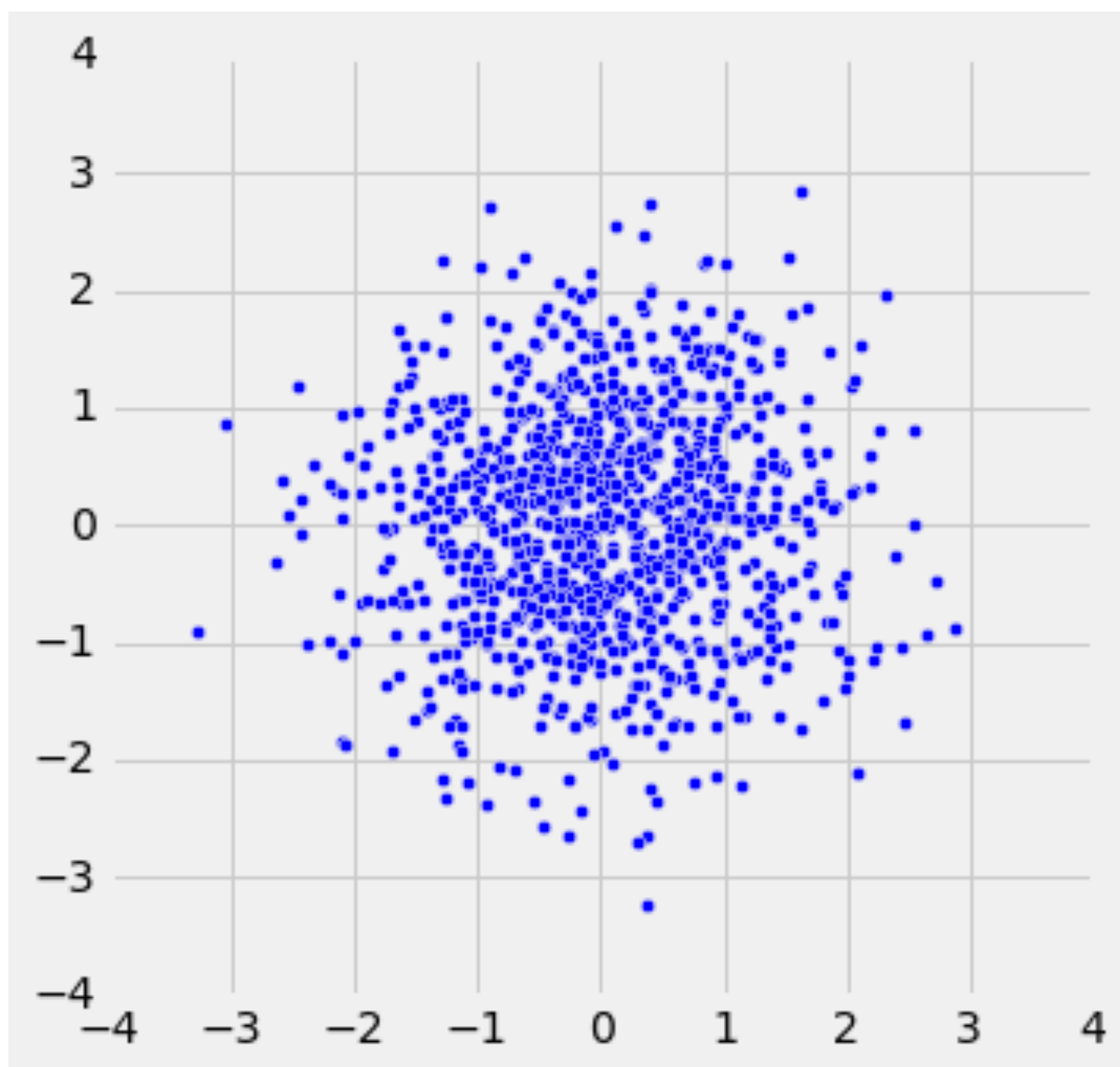
Therefore the $\text{Pearsonr}(x,y)$ is just the linear regression slope when the variables are rescaled to have standard deviation equal to 1.



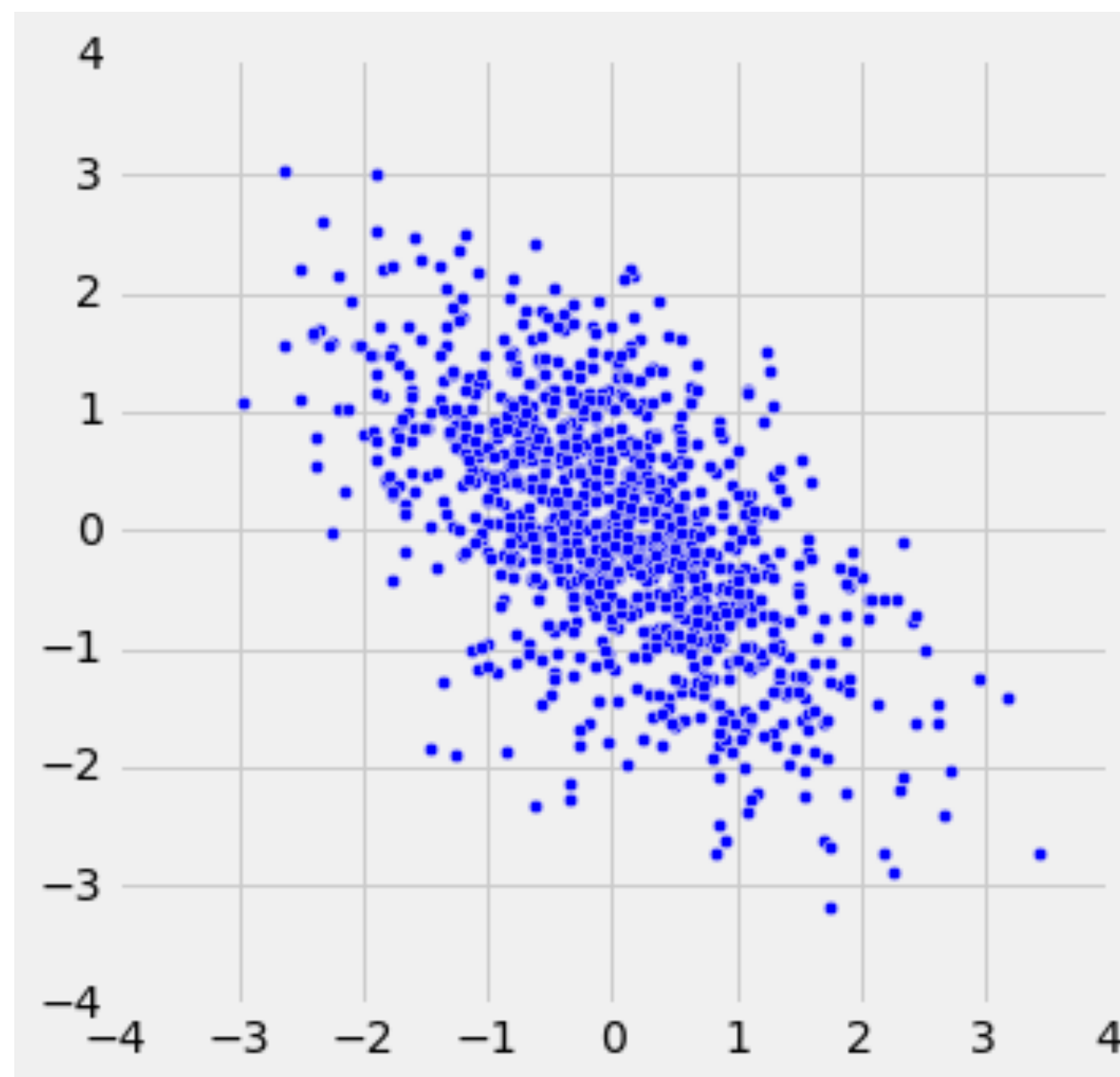
.9



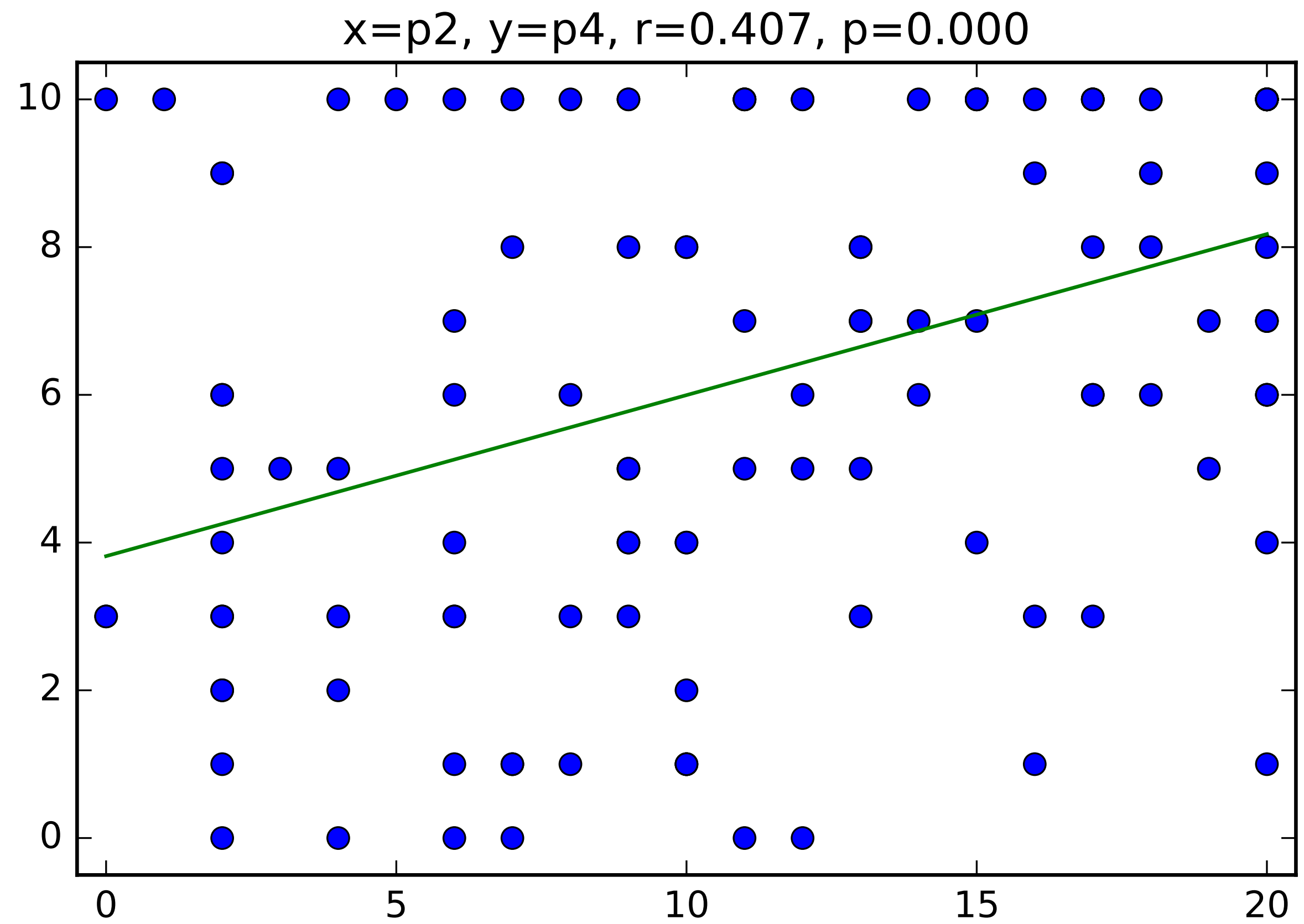
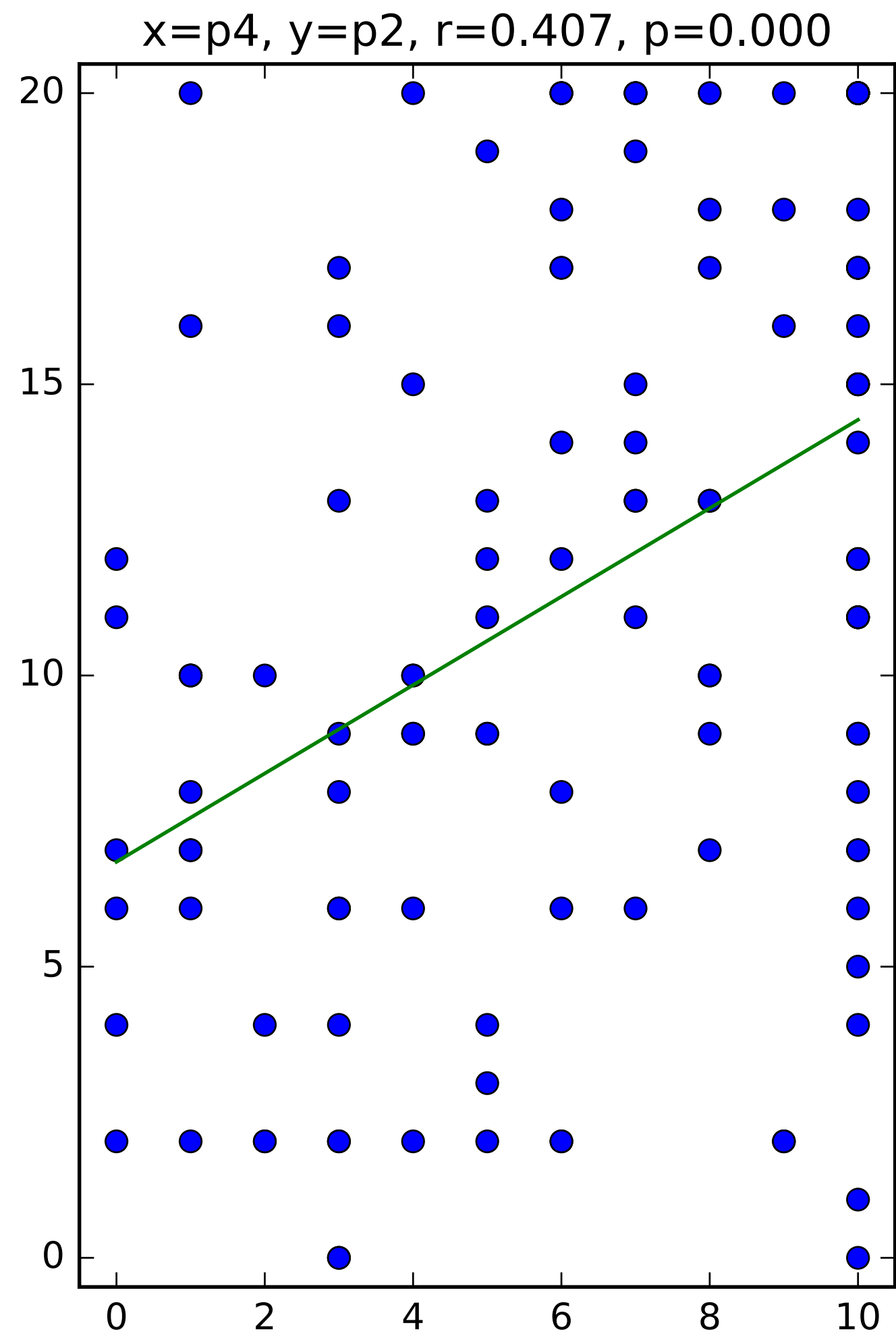
.25



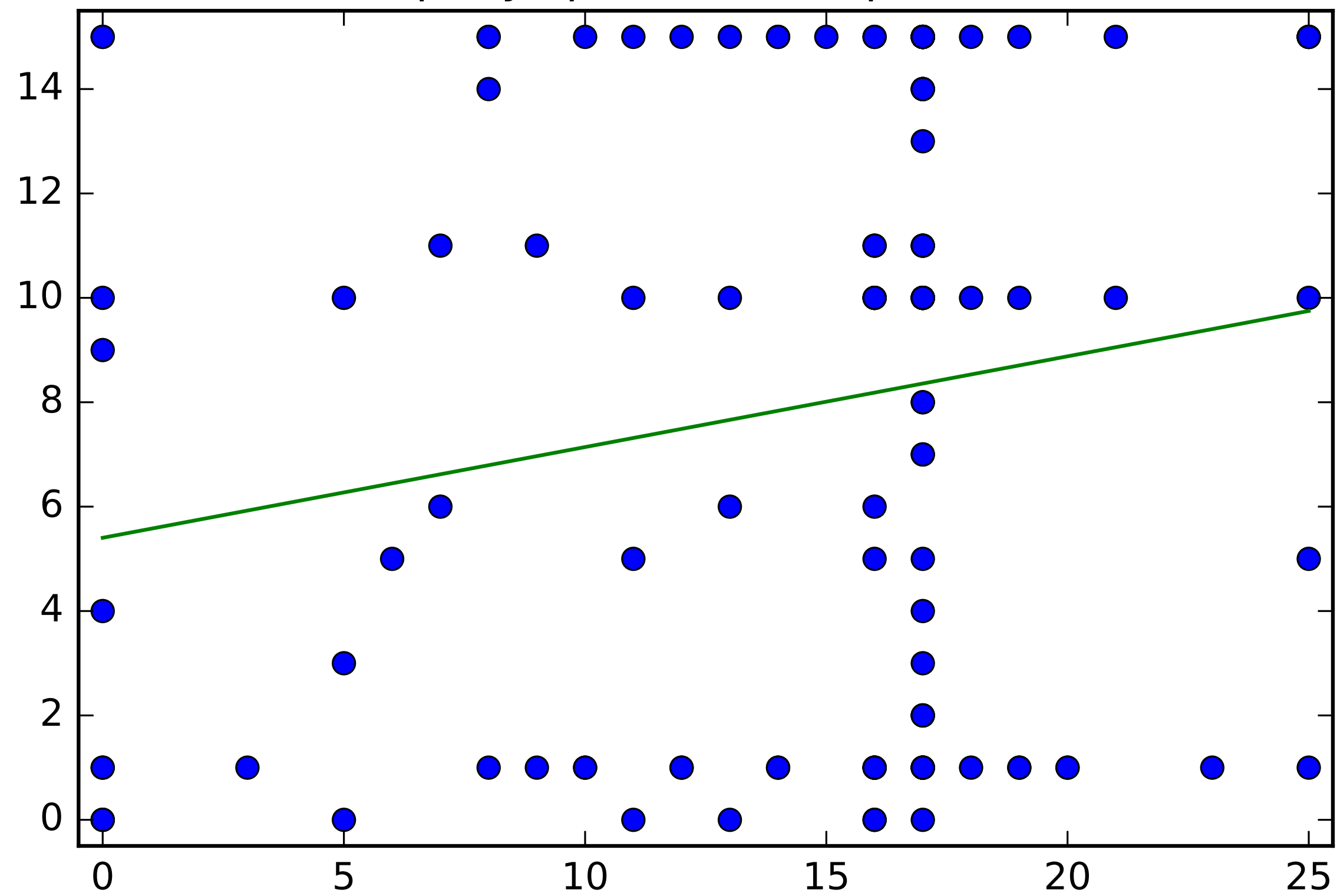
0

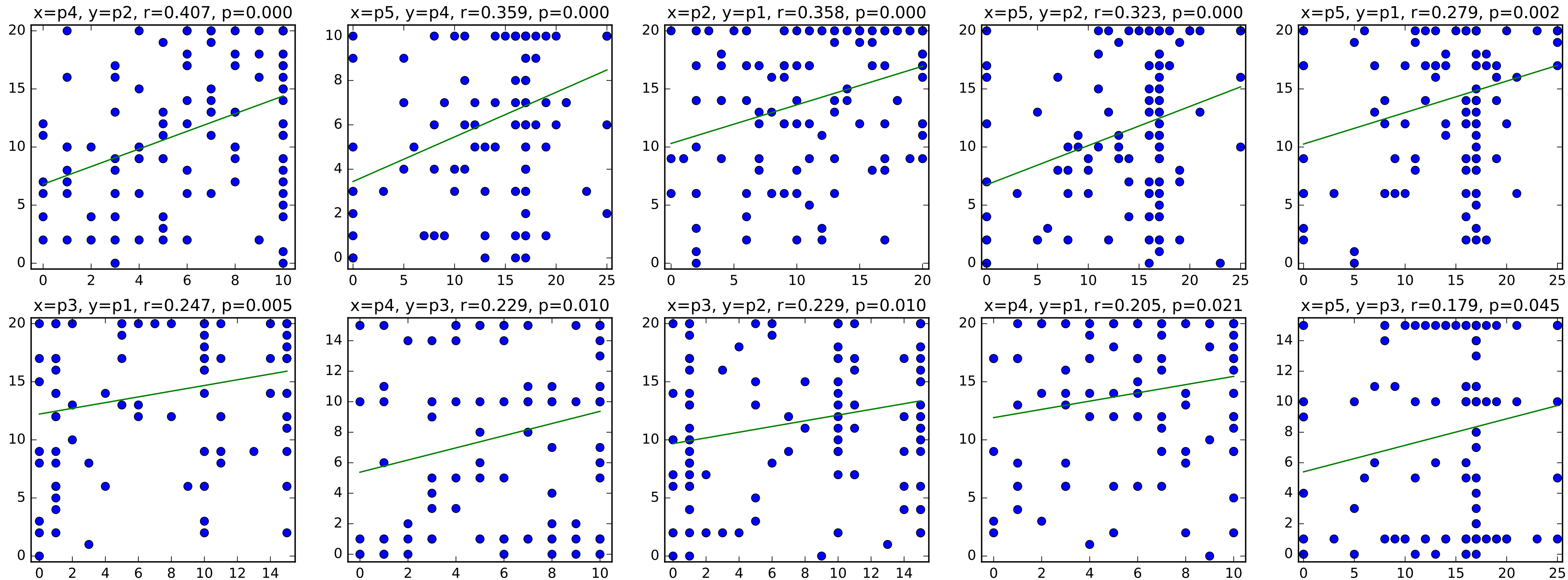


-.55



$x=p5, y=p3, r=0.179, p=0.045$

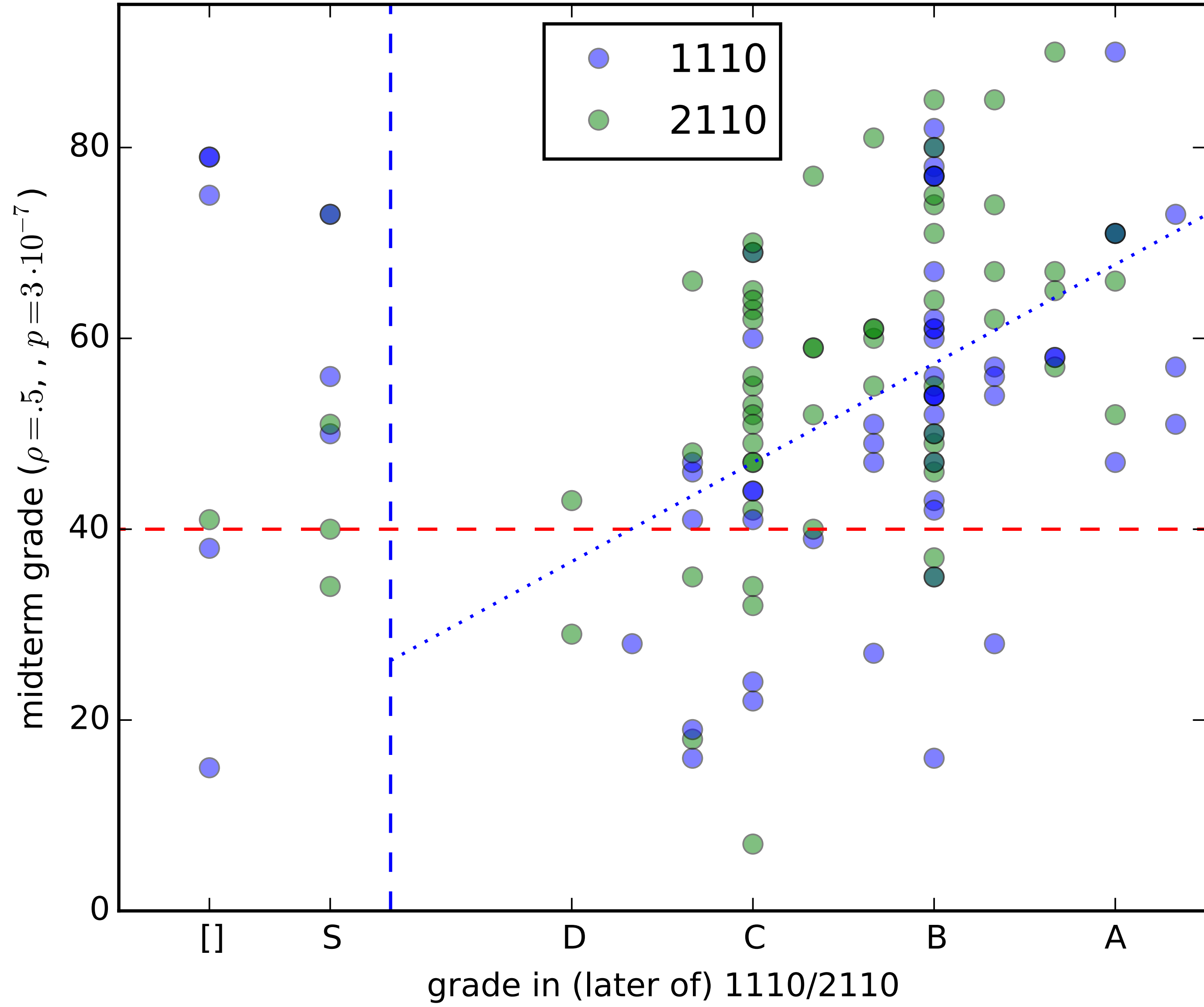




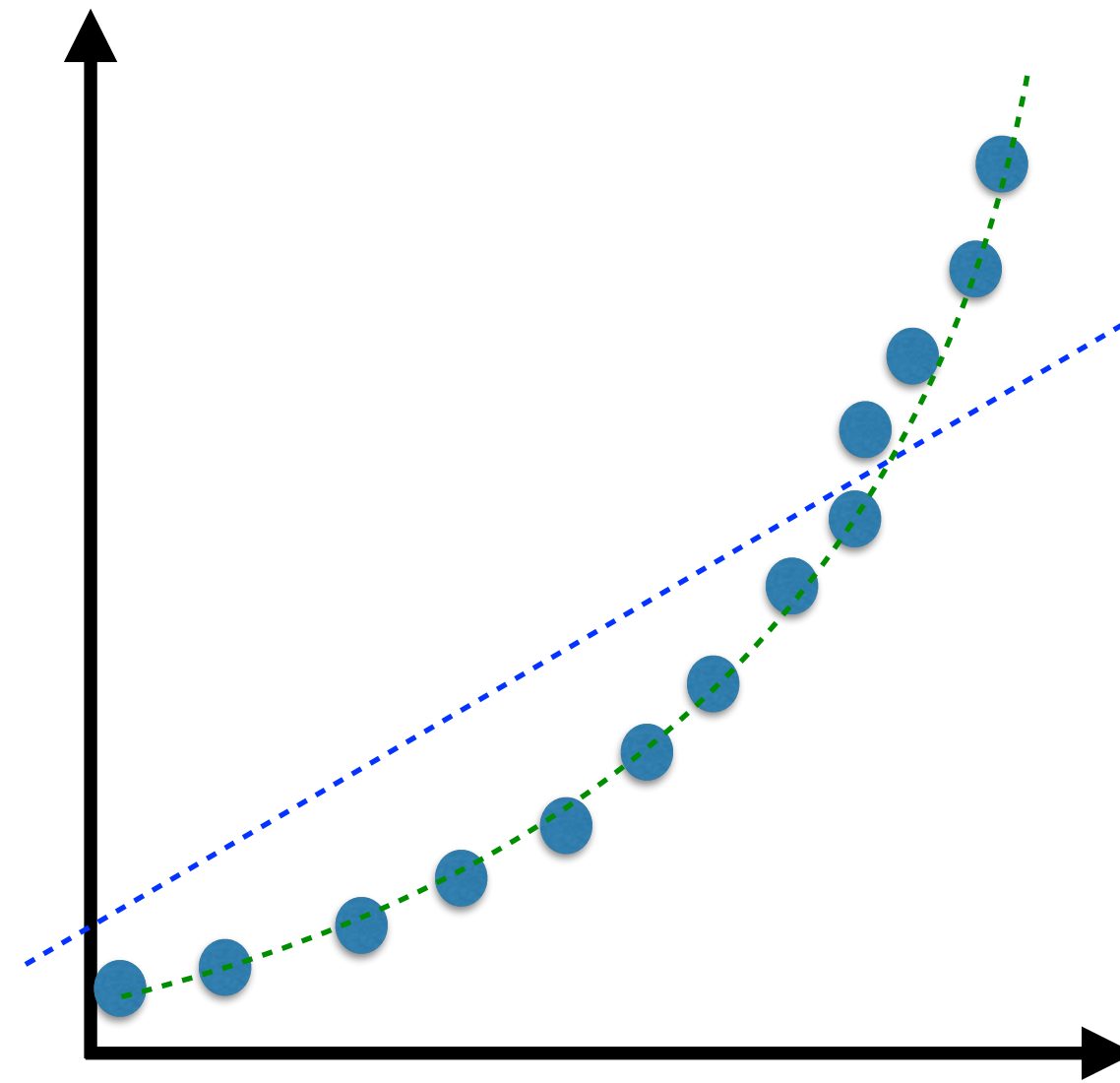
[20,20,15,10,25]

1) naive bayes, 2) simple probab, 3) wikipedia, 4) python, 5) research lit

students in info2950 spr'17



Dependences are not always linear



(Except in ...)

**GENERAL ASSEMBLY OF NORTH CAROLINA
SESSION 2011**

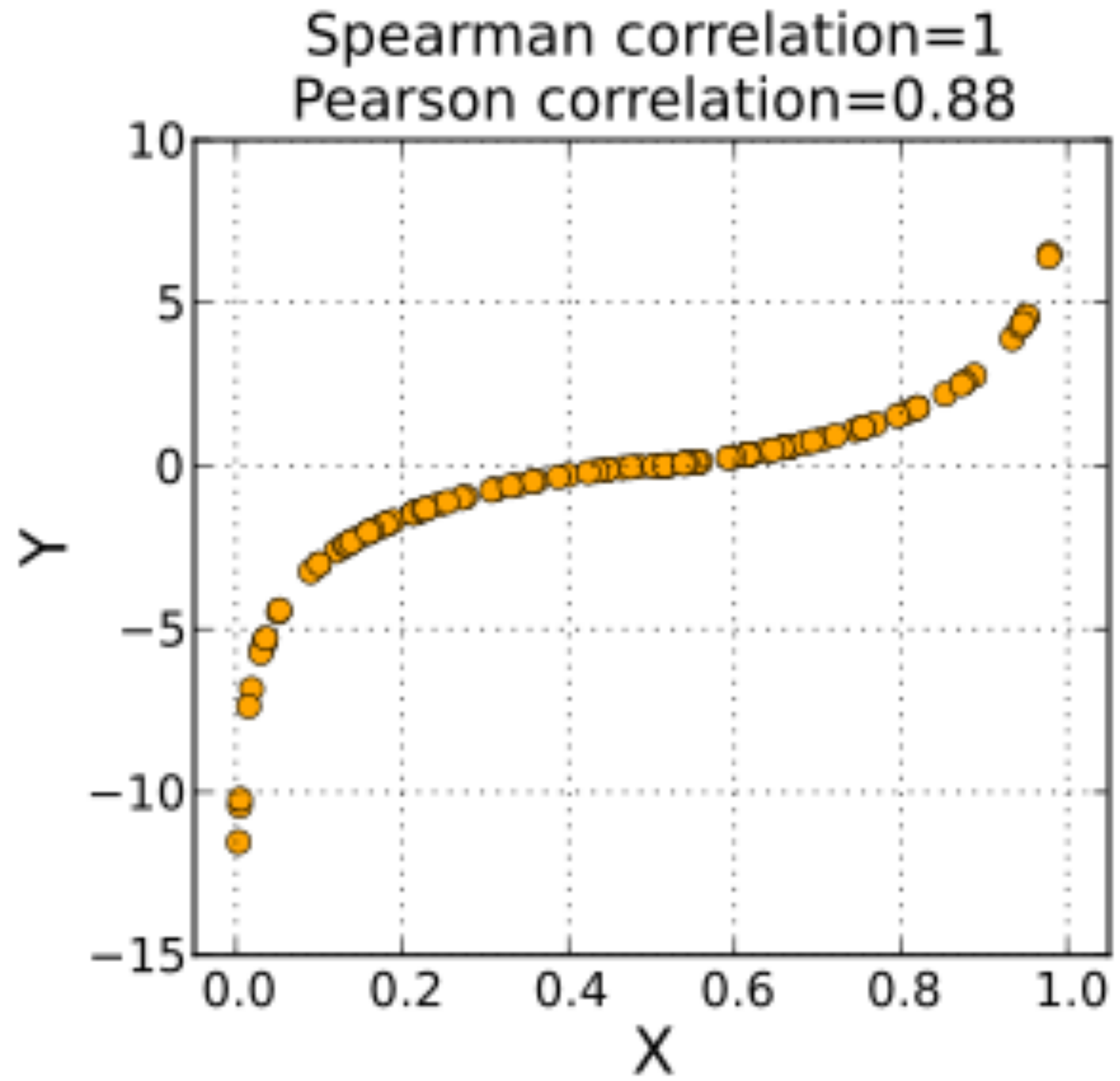
H

D

**HOUSE BILL 819
Committee Substitute Favorable 4/21/11
PROPOSED SENATE COMMITTEE SUBSTITUTE H819-CSLH-38 [v.18]**

4/25/2012 11:05:28 AM

10 (e) The Division of Coastal Management shall be the only State agency authorized to
11 develop rates of sea-level rise and shall do so only at the request of the Commission. These
12 rates shall only be determined using historical data, and these data shall be limited to the time
13 period following the year 1900. Rates of sea-level rise may be extrapolated linearly to estimate
14 future rates of rise but shall not include scenarios of accelerated rates of sea-level rise. Rates of
15 sea-level rise shall not be one rate for the entire coast but, rather, the Division shall consider
16 separately oceanfront and estuarine shorelines. For oceanfront shorelines, the Division shall use
17 no fewer than the four regions defined in the April 2011 report entitled "North Carolina Beach



The Pearson correlation coefficient misses non-linear relationships and is also sensitive to outliers — the Spearman correlation can sometimes find correlations that Pearson misses.

It is defined as the Pearson correlation of the rank order of the data.

That means it also varies from -1 (perfectly anti-correlated) to $+1$ (perfectly correlated), with 0 meaning uncorrelated.

If the data has $x = [.6, .4, .2, .1, .5]$ then the ranks are $r = [5, 3, 2, 1, 4]$.

For data $y = [403, 54, 7, 2, 148]$, the ranks $s = [5, 3, 2, 1, 4]$ are the same.*

so the Spearman correlation is 1, whereas the Pearson is less than one.

Both functions are available in `scipy.stats` (as `pearsonr()` and `spearmanr()`).

[*Actually the second was generated from the first by taking the integer part of $\exp(10x)$]

Defined as the Pearson correlation for the ranks, the Spearman correlation is written

$$\rho = \frac{\text{Cov}[r, s]}{\sigma[r]\sigma[s]}, \quad (1)$$

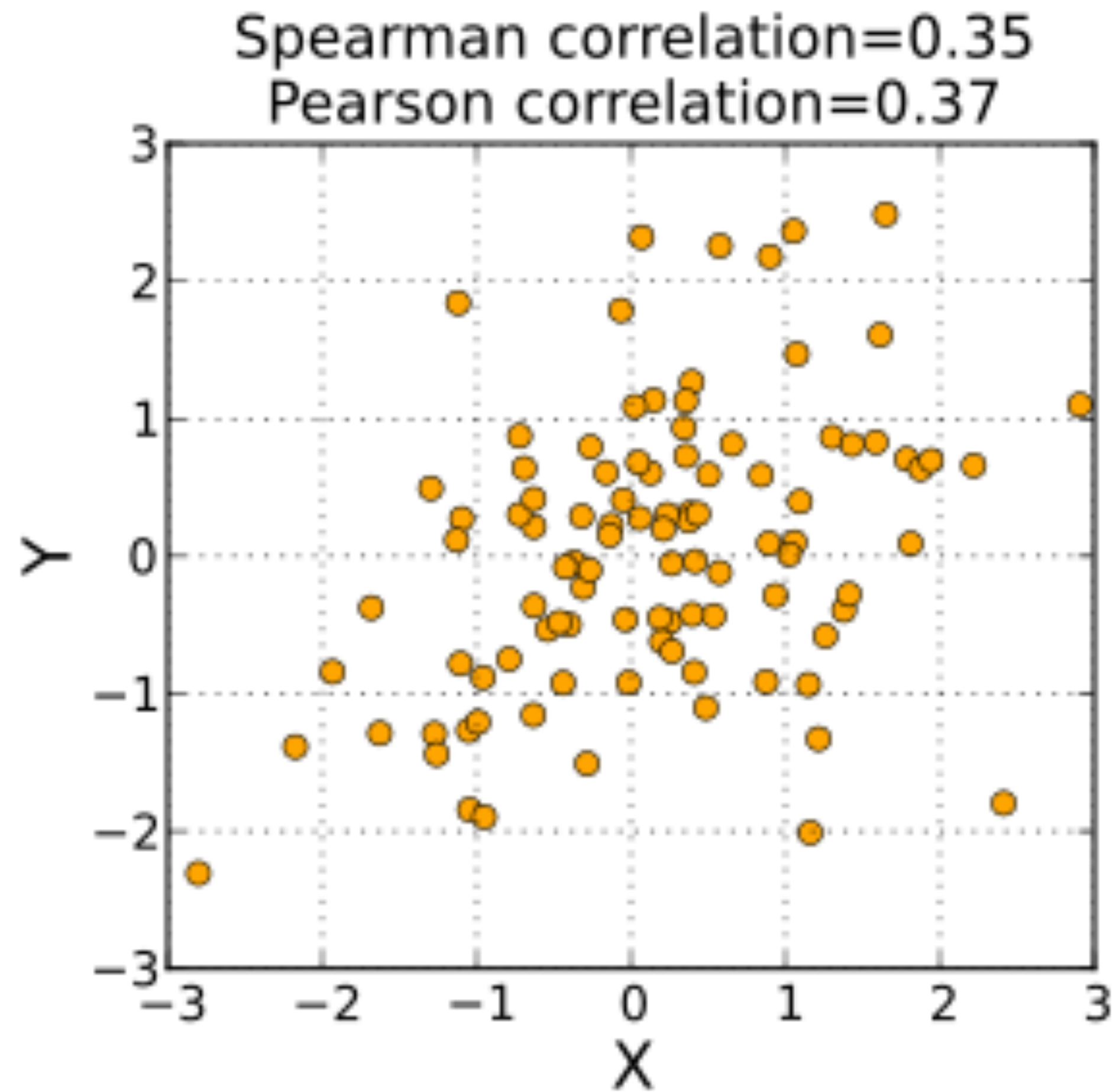
where $\text{Cov}[r, s] = E[(r - E[r])(s - E[s])]$

(generalizing the $\text{Var}[x] = E[(x - E[x])^2]$, with $\text{Cov}[x, x] = \text{Var}[x]$).

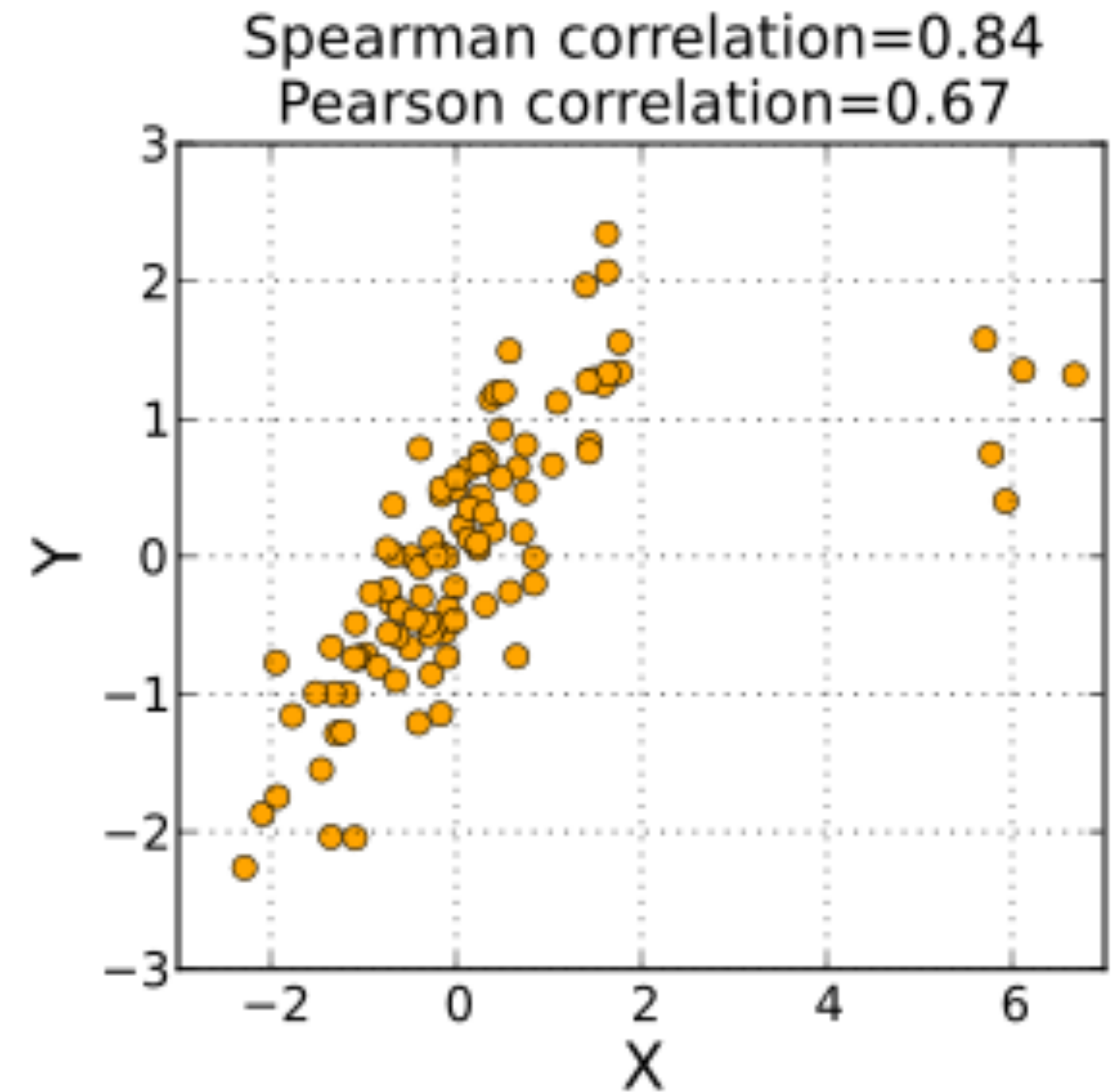
The formula for the Spearman correlation coefficient is given at http://en.wikipedia.org/wiki/Spearman's_rank_correlation_coefficient in terms of the difference $d_i = r_i - s_i$ between ranks, in this easily calculable form:

$$\rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}. \quad (2)$$

It is straightforward to verify that (1) reduces to (2) (see linked notes)



roughly elliptically
distributed and there are
no prominent outliers:
same



less sensitive than the Pearson
correlation to strong outliers that are
in the tails of both samples
(limited to value of rank)

The Democrats Show Some Spunk

Two days of obstruction on Capitol Hill.

By *Jim Newell*

We are on Day 12 of the Trump administration, and Senate Democrats have begun slowing the pace of Capitol Hill down to a crawl.



It started Monday night with a surprise procedural move, just as Congress was returning to session. The Senate Finance Committee was scheduled to meet at 4 p.m. to



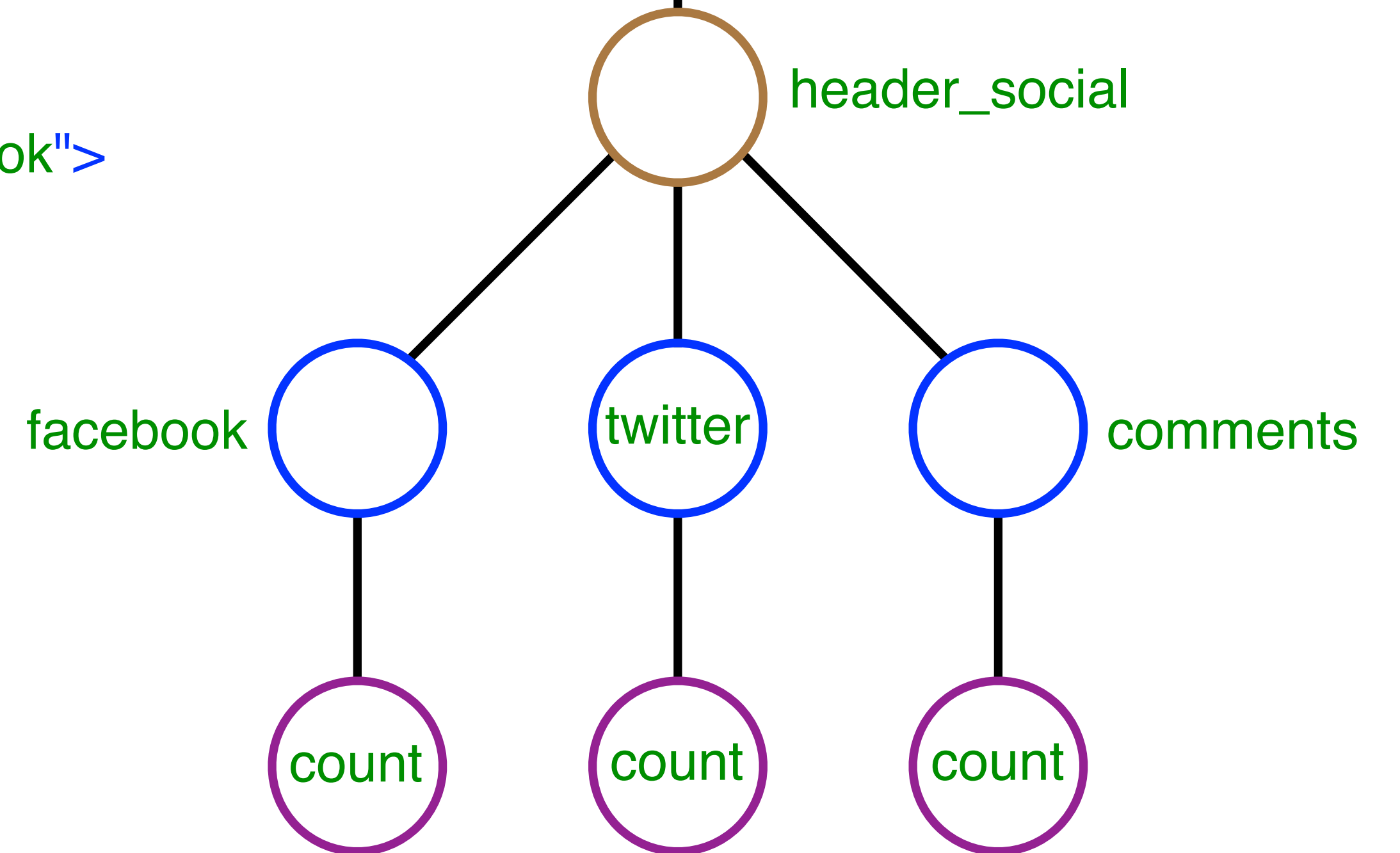
dynamically generated, so scrape with selenium

```
<div class="social social-with-popup sharing-buttons header_social">
```

```
<a href="#" data-share="http://www.facebook.com/..." class="facebook">  
  <span class="icon"></span>  
  <span class="count">8.9k</span>  
</a>
```

```
<div class="twitter">  
  <a href="#" data-share="http://twitter.com/..." class="twitter">  
    <span class="icon"></span>  
  </a>  
  <a class="tw-count" href="http://twitter.com/search?q=...">  
    <span class="count">2.1k</span>  
  </a>  
</div>
```

```
<a href="#comments" class="comments">  
  <span class="icon"></span>  
  <span class="count">1.4k</span>  
</a>  
</div>
```



```
allsoc=[]  
for url in urls:  
  driver.get(url)  
  social = driver.find_element_by_class_name('header_social')  
  soc=[] #temporary list for this url  
  for cl in ('facebook','twitter','comments'):  
    elt = social.find_element_by_class_name(cl)  
    subelt = elt.find_element_by_class_name('count')  
    soc.append(subelt.get_attribute('innerHTML'))  
  allsoc.append(soc)
```