

Here is a proof that the correlation between X and Y' is 1.

1. To simplify ... assume that  $x = (X - \text{mean } X)$ ,  $y = (Y - \text{mean } Y)$

2. General formula for  $r_{xy} = \hat{\sigma}_{xy} / N * \sigma_x * \sigma_y$

3. Now, in z score form, we would have ...  $\hat{\sigma}_{z_x * z_y} / N$

4. What if what we want is the correlation between X and Y' .... so

5. The formula would look like ....  $r_{XY'}$  ... or  $r_{z_x z_y}$

6. Keep in mind that ...

$$Z_{y'} = r_{xy} * Z_x$$

and

$$\sigma_{z_x} = \sigma_{z_y} = 1$$

7. Now ... if the standard deviation of  $Z_x$  is 1, then the standard deviation of  $r_{xy} * Z_x$  would be like having a variable ( $Z_x$ ) multiplied by a constant of  $r$  ... and whenever you multiply a variable by a constant ... the standard deviation is impacted by the value of the constant. So ... since  $Z_{y'} = r_{xy} * Z_x$  .... then whatever happens to the standard deviation of  $Z_x$  would also happen to  $Z_{y'}$ . If the constant is  $r$  ... then the standard deviation of  $Z_{y'}$  would also be impacted by  $r$  ... and rather than being 1 if it were  $Z_z$  or  $Z_y$  ... would only be the size of  $r * 1$  .... or only  $r$ .

This is a long way of saying that  $\sigma_{y'} = r_{xy}$

8. At long last, we can now work on the  $r$  formula between X and Y'.

$r_{z_x z_y} = \hat{\sigma}_{z_x * z_y} / N$  .... but, in the denominator ... while the standard deviation of  $z_x$  is 1 ... and we don't need it, the standard deviation of  $z_y$  is not 1 ... but rather only  $r_{xy}$ . Thus, what we really have as the formula

is ....

$$r_{xy} = \frac{\sum z_x * z_y}{N * r_{xy}}$$

9. Continuing ... recall from 6, we could substitute  $r_{xy} * z_x$  for  $z_y$  and we would have ...

$$r_{z_x z_y} = \frac{\sum z_x * r_{xy} * z_x}{N * r_{xy}}$$

10. This would lead to a simplification

$$= \frac{\sum z_x * z_x * r_{xy}}{N * r_{xy}}$$

11. Since the rs cancel out from top to bottom ... we would have left

$$r_{z_x z_y} = \frac{\sum z_x z_x}{N}$$

12. Since z scores are deviation scores around the mean of 0, this means that  $\sum z_x^2 / N$  is merely the average of the squared deviations using z score data ... and that is the variance ... and that = 1!

13. Finally! ....  $r_{xy} = 1$

And that is why you got the scatterplot between predictor and predictED of a straight line ...