

Prep-assessment poll #1: In the formula for r-squared, what is the meaning of the ratio between sum of squared residuals (SSR) and sum of squared differences (SST) from mean?

Answer #1: Let's compare cases when the regression line for the data is parallel to Ox. Then, the SST would be a relatively small number (compared to when data is on a line with some slope), meaning that any residuals in SSE would be weighted more heavily, since we're dividing by a smaller number. Then, our  $R^2$  coefficient would be "worse" (we need more precise fit). When data is on a slope, we would like to account for the fact that residuals are increased due to the slope (since we're not doing perpendicular projections to the line of best fit), so the term SST will be larger, making  $R^2$  "better."

Prep-assessment poll #2: How do we find the linear line of best fit to the data?

Answer #2: To do so, we need to minimize mean squared error (MSE) or sum of squared residuals (SSR). This error function is the sum of squares of difference between actual  $y_0$  at point  $x_0$  and predicted  $\hat{y}_0$  (where  $\hat{y}(x) = ax + b$ ) at  $x_0$ . We then would like to differentiate this function w.r.t.  $a$  and  $b$  and equalize those to 0 (then, this would mean that we're in a local minimum). Solving the equations, we get the coefficients that would minimize the MSE and be the linear regression line for our data.