

Principal Component Analysis (PCA)

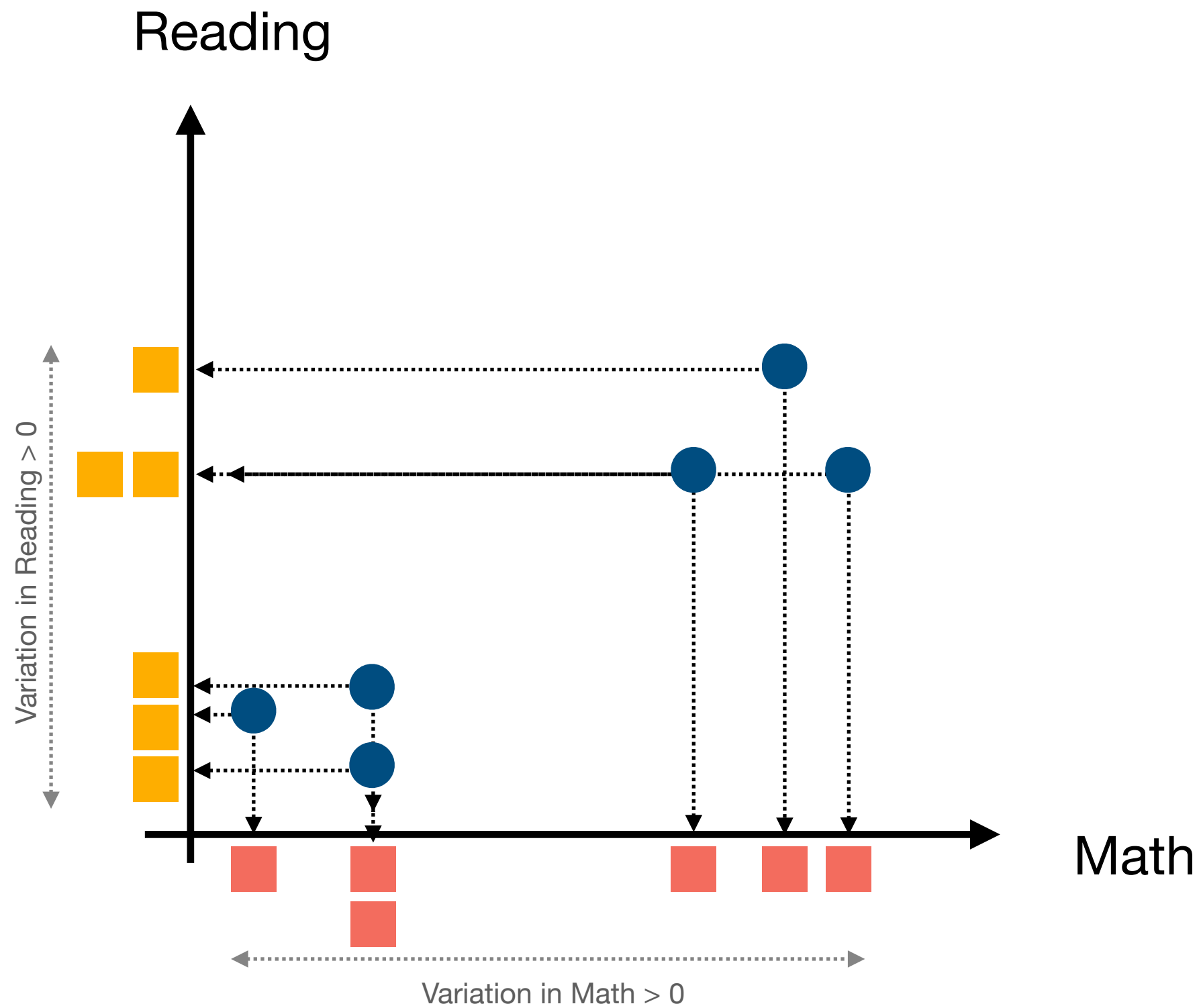
- Dimension Reduction
- Representation of Data

อ.ดร.สัวะโชติ ศรีสุทธียากร และ อ.ดร.ประภาศิริ รัชชประภาพรกุล
ภาควิชาวิจัยและจิตวิทยาการศึกษา คณะครุศาสตร์
จุฬาลงกรณ์มหาวิทยาลัย

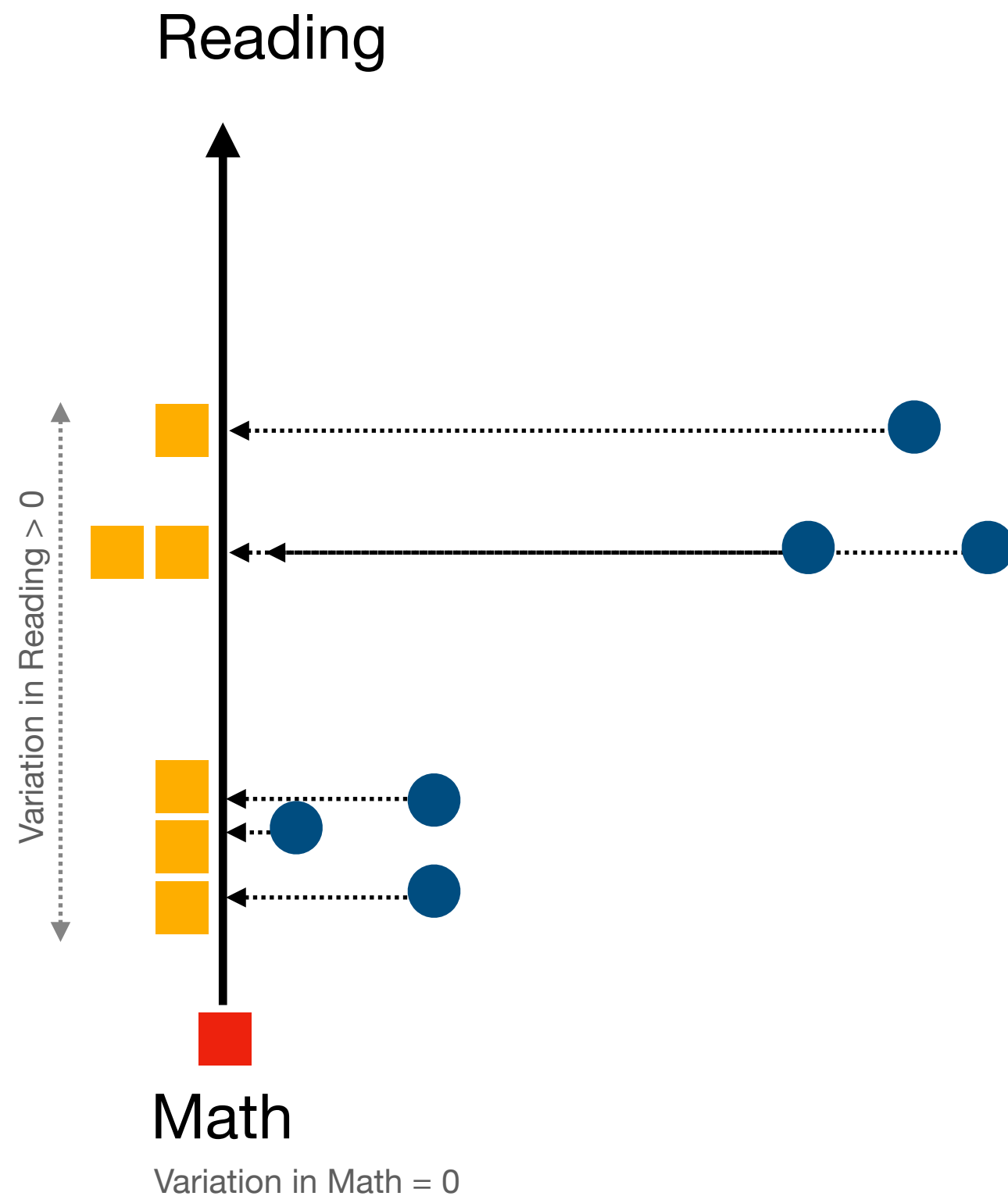
PCA - concept

Student	Reading	Math
บุญมี	8	99
บุญมา	5	43
บุญหนัก	1	16
บุญทับ	4	49
บุญถึก	7	83
บุญถึง	3	55

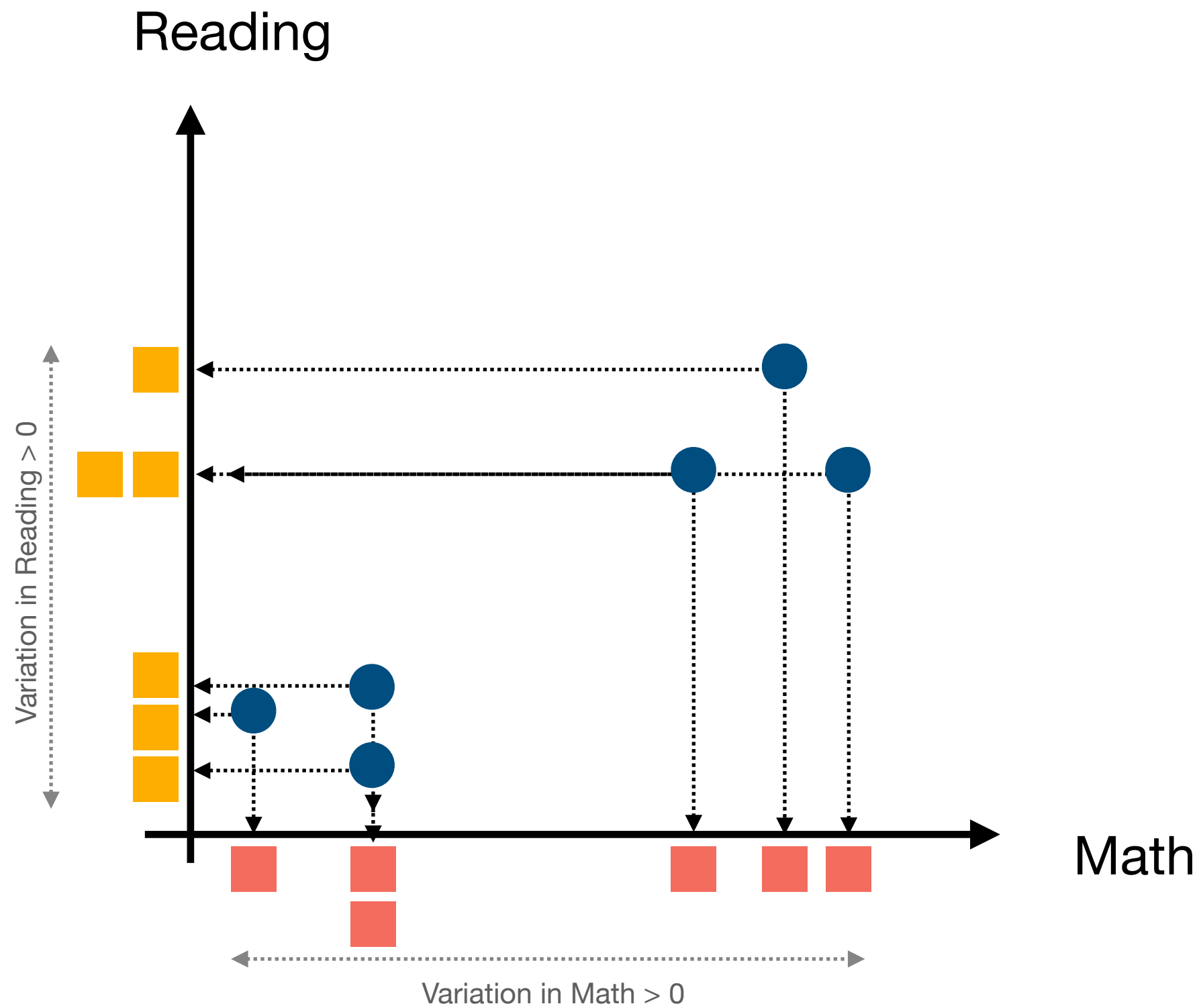
PCA - concept



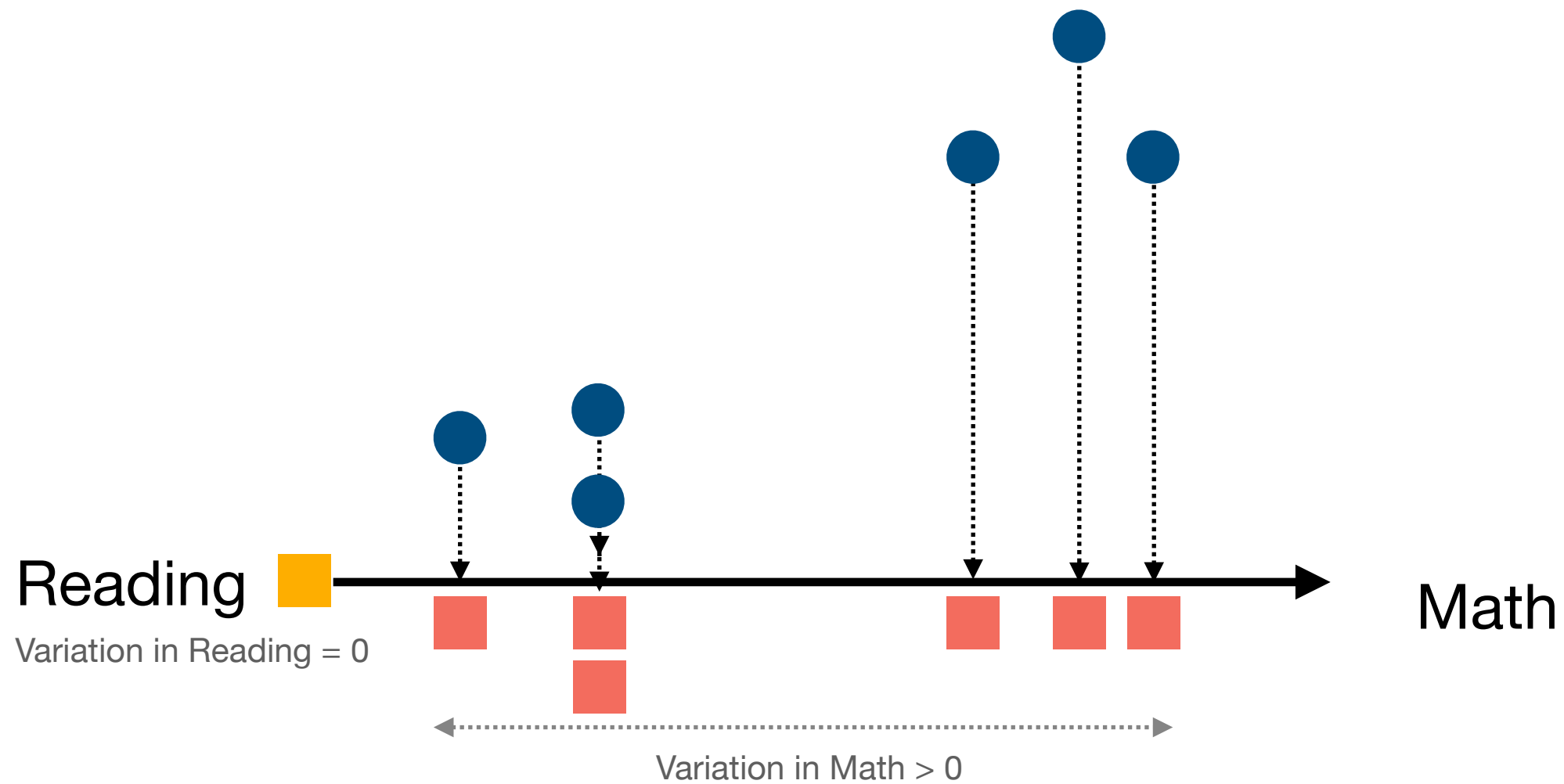
PCA - concept



PCA - concept



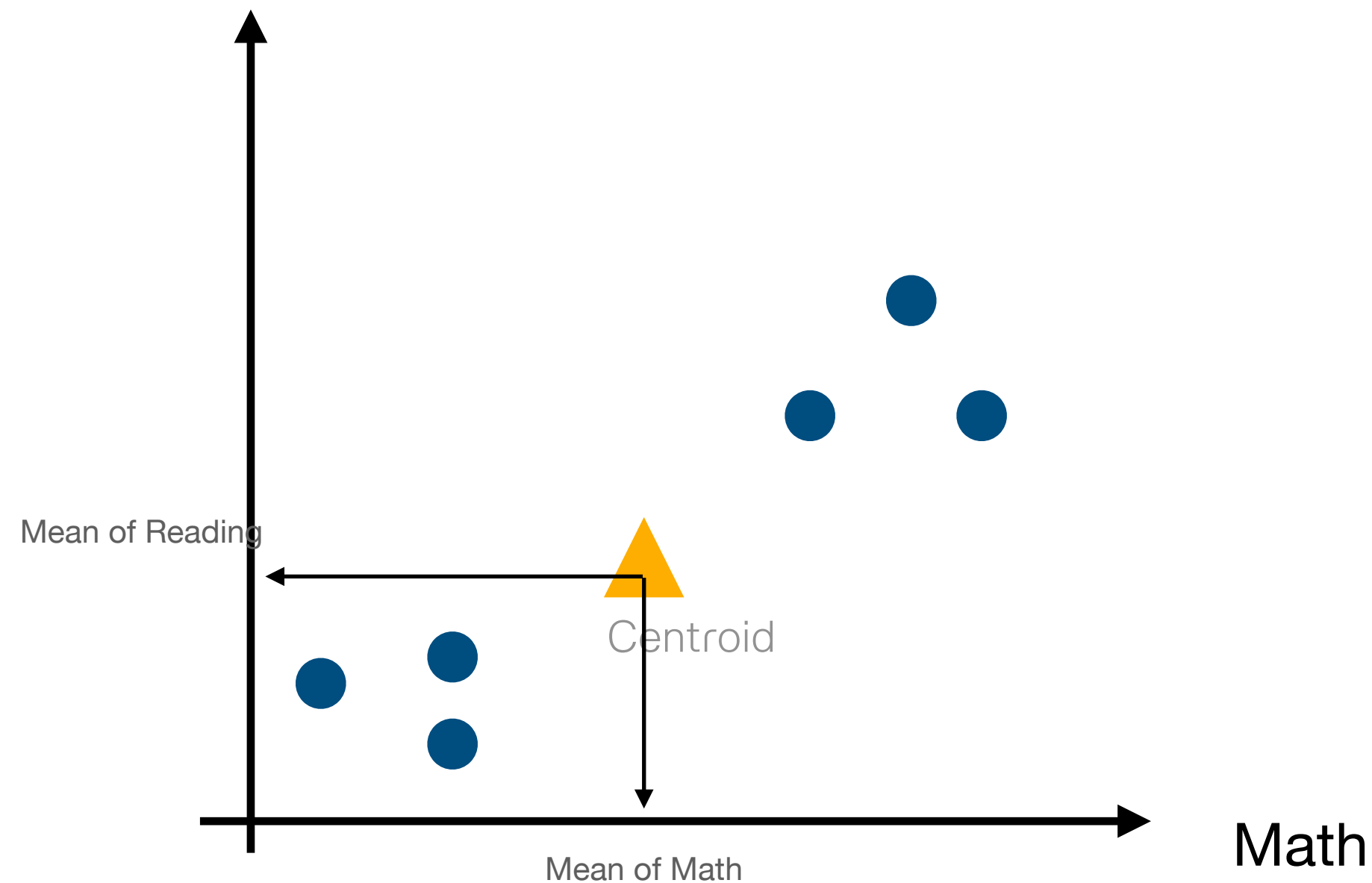
PCA - concept



PCA - concept

Calculate centroids and standard deviations of the data

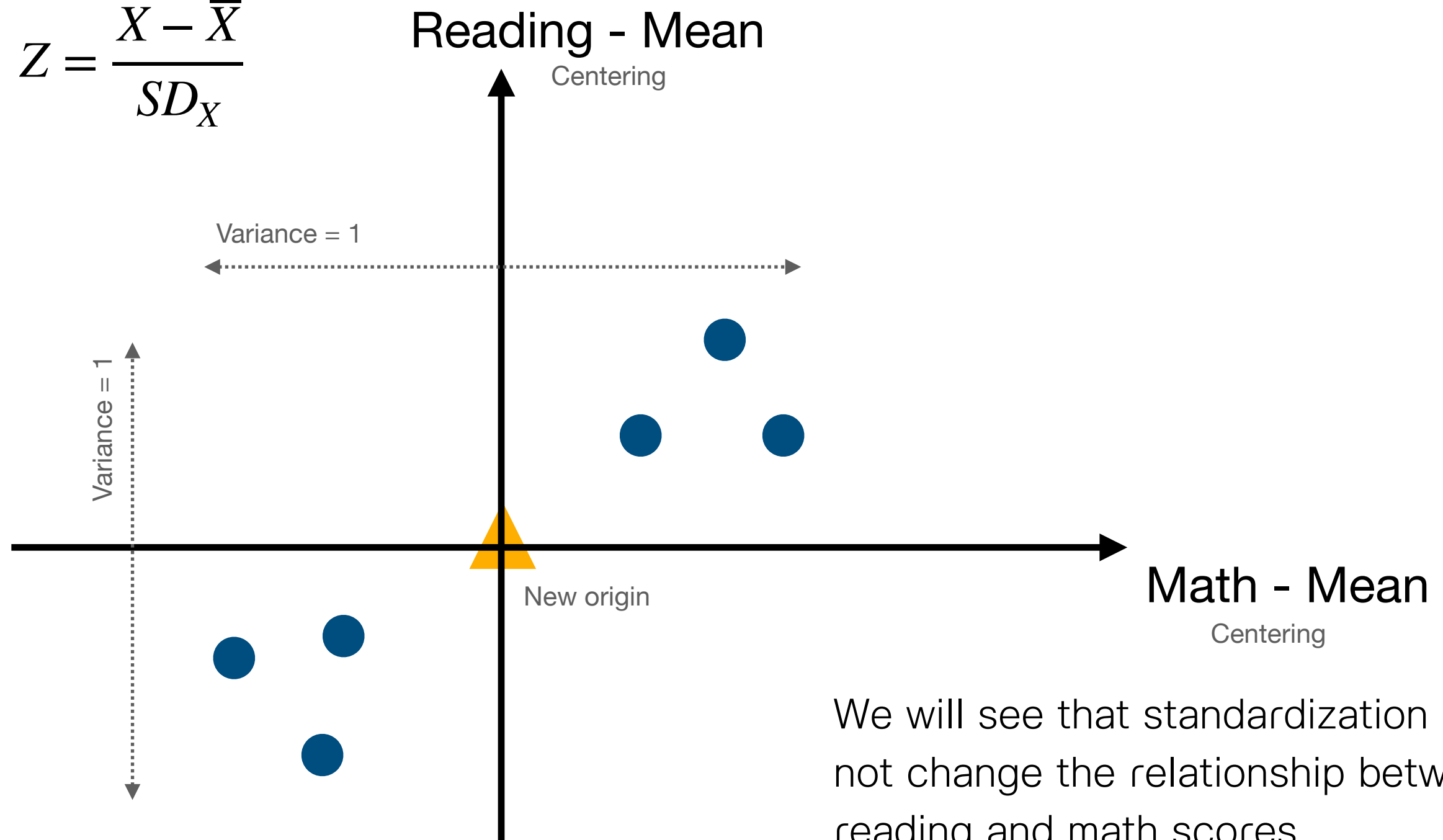
Reading



PCA - concept

Standardized the variables

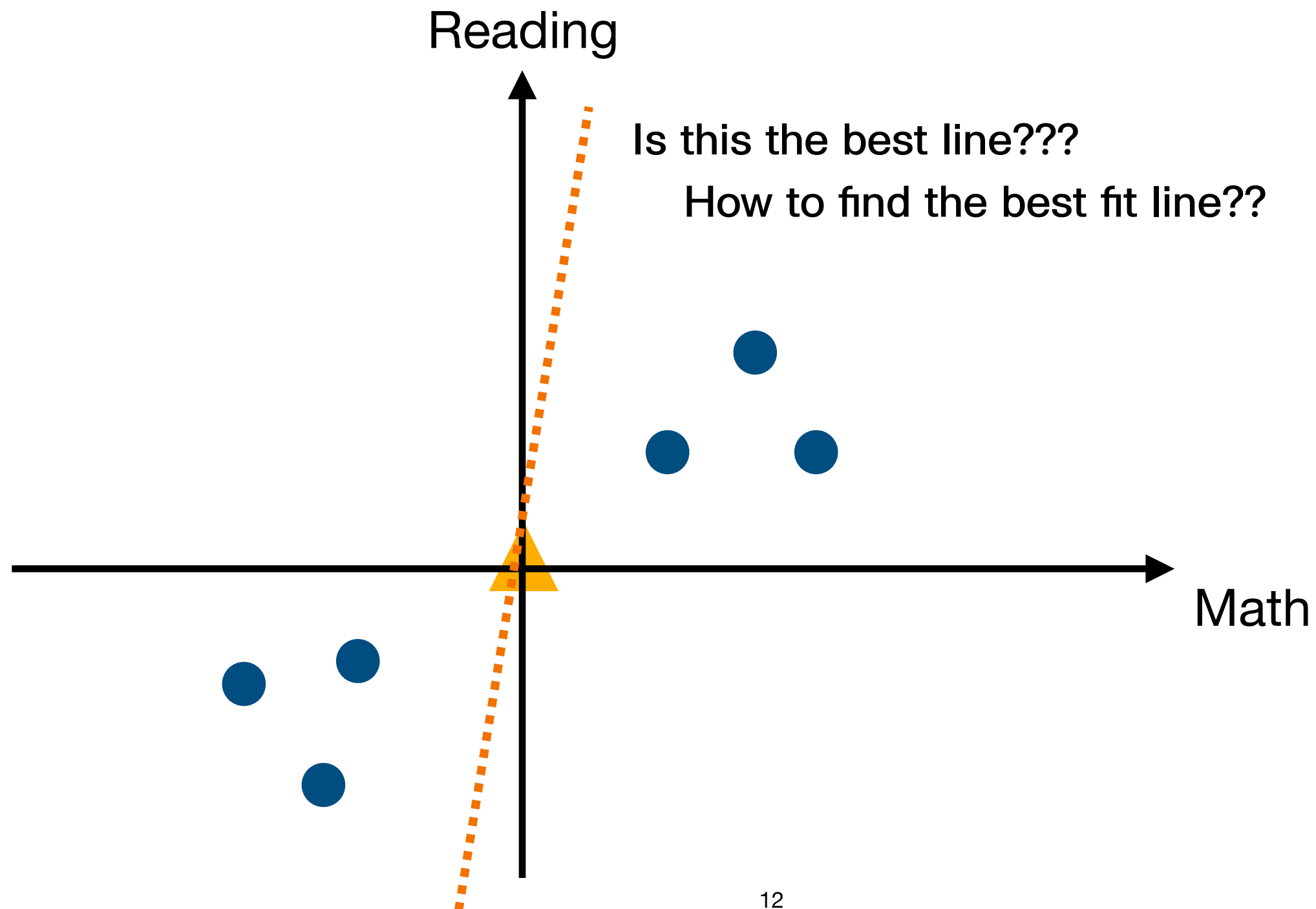
$$Z = \frac{X - \bar{X}}{SD_X}$$



We will see that standardization did not change the relationship between reading and math scores.

PCA - concept

Fit a line to the data...

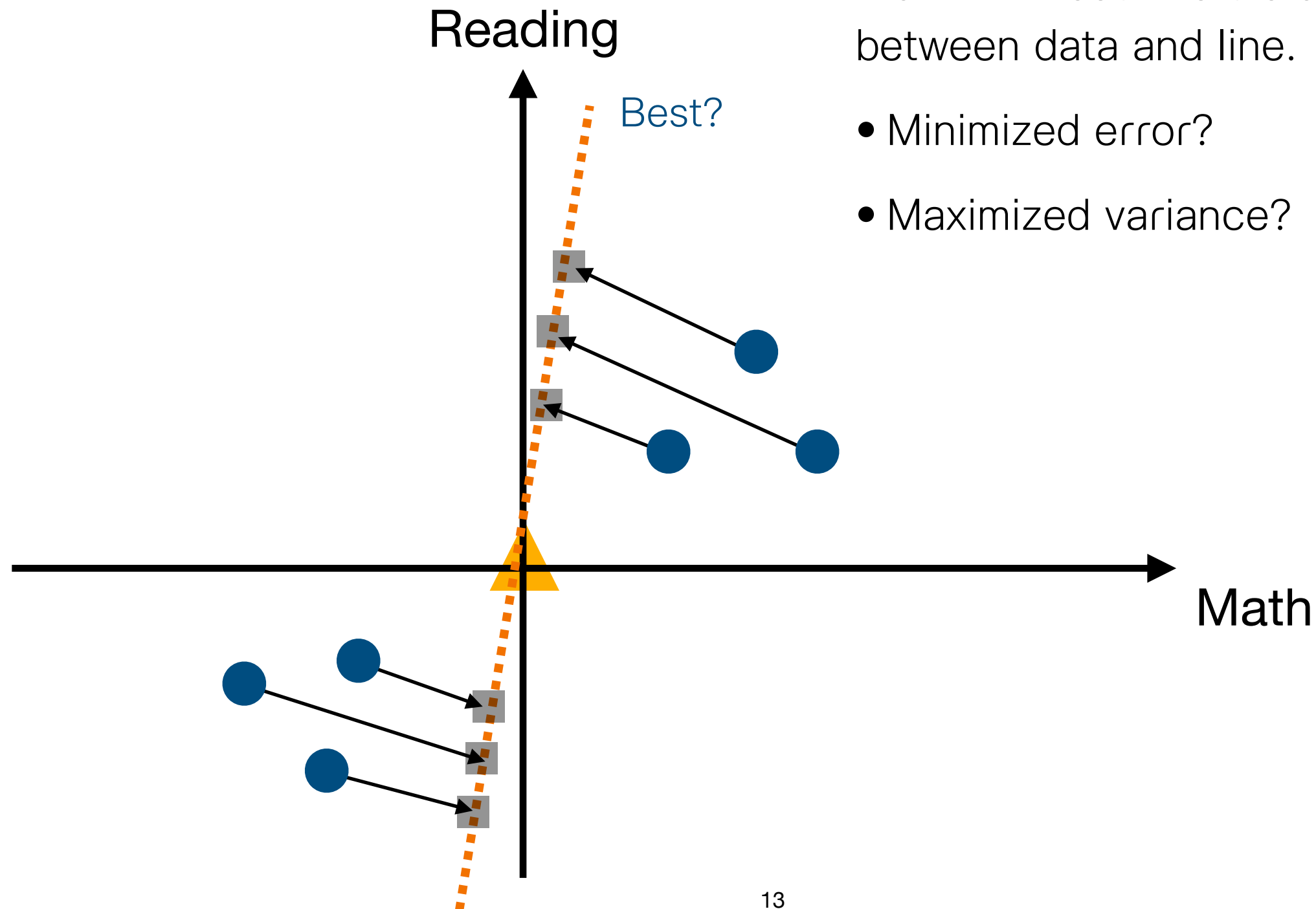


PCA - concept

Fit a line to the data...

The best fit line is the line that minimised the distances between data and line.

- Minimized error?
- Maximized variance?

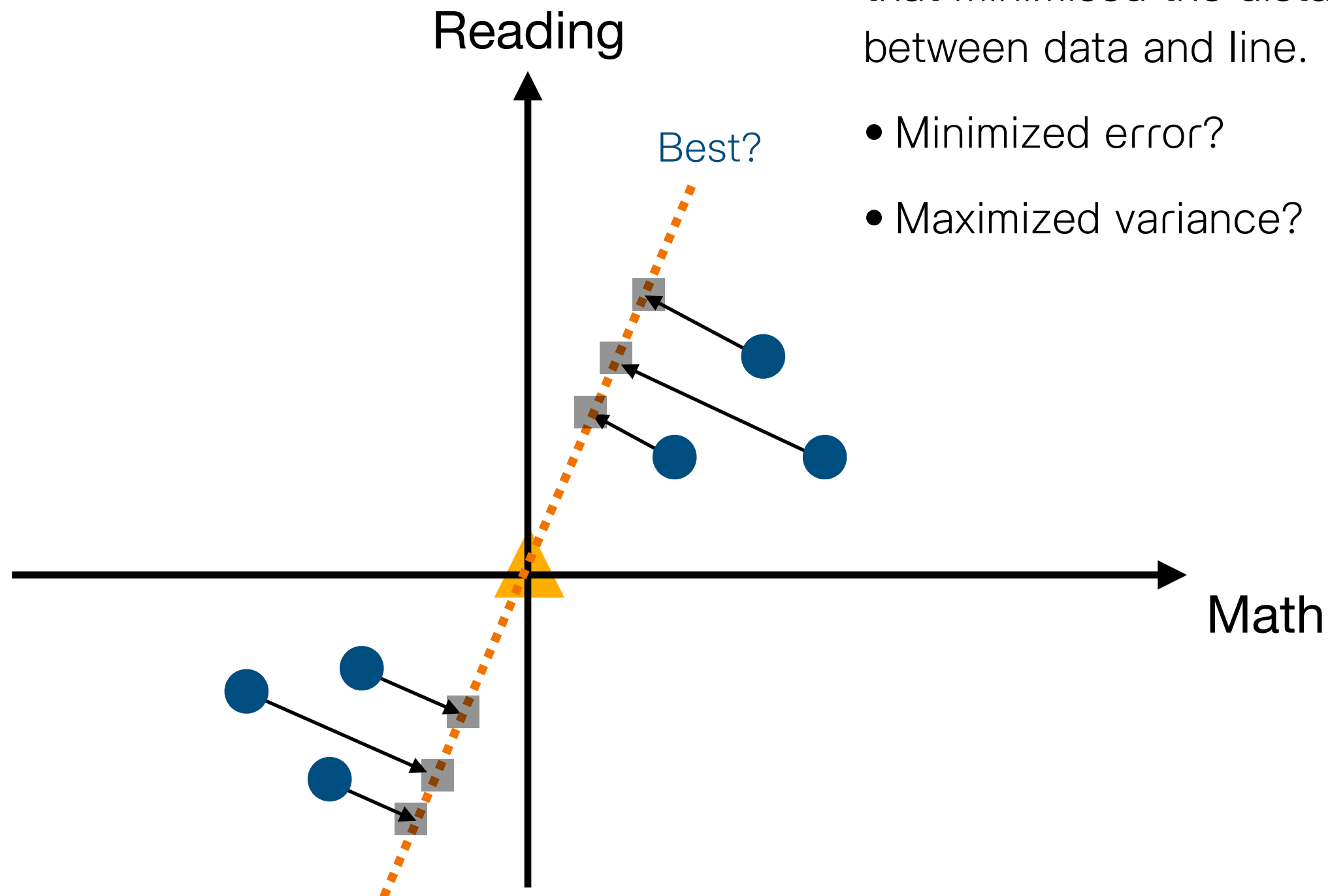


PCA - concept

Fit a line to the data...

The best fit line is the line that minimised the distances between data and line.

- Minimized error?
- Maximized variance?

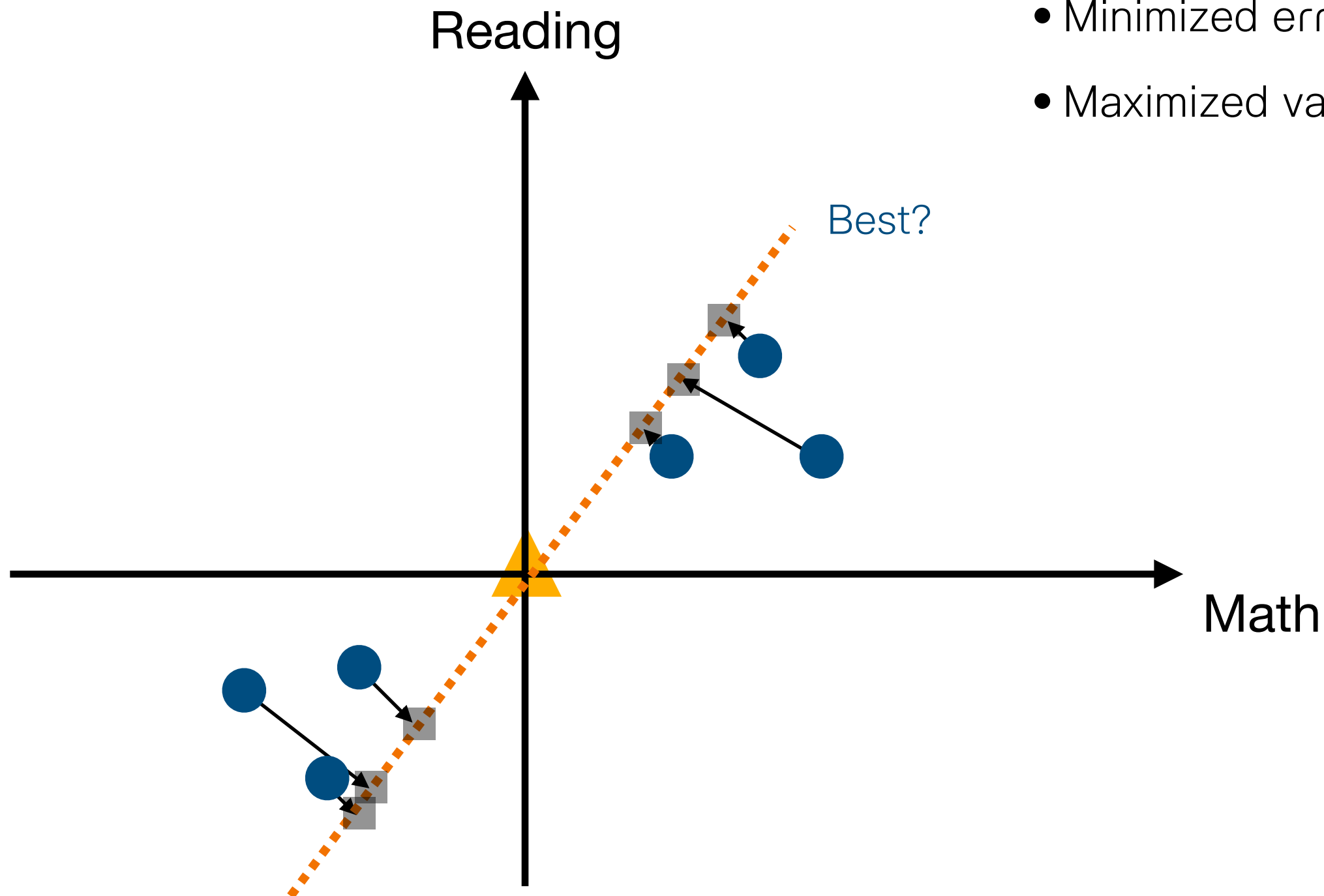


PCA - concept

Fit a line to the data...

The best fit line is the line that minimised the distances between data and line.

- Minimized error?
- Maximized variance?

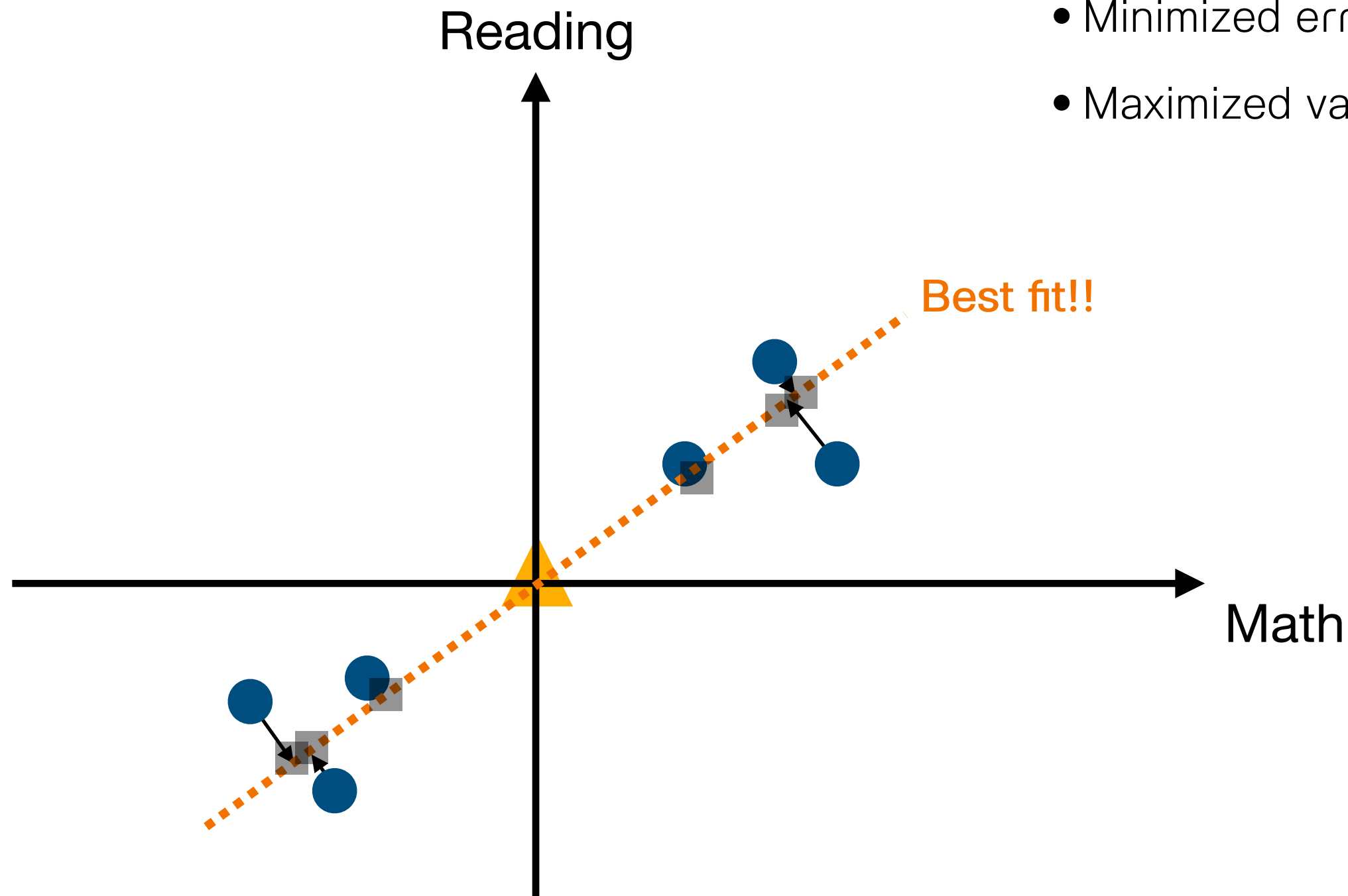


PCA - concept

Fit a line to the data...

The best fit line is the line that minimised the distances between data and line.

- Minimized error?
- Maximized variance?

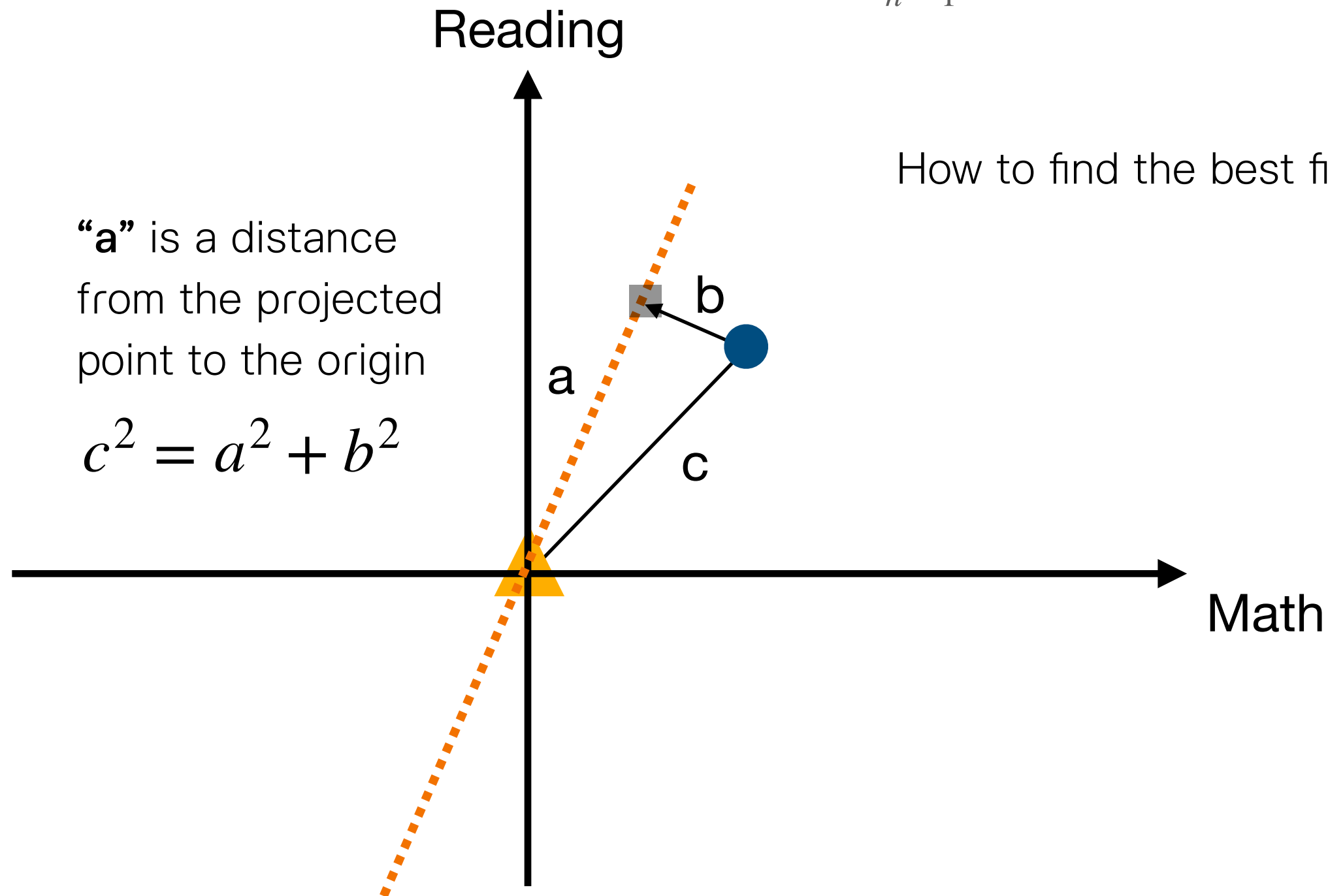


PCA - concept

Fit a line to the data...

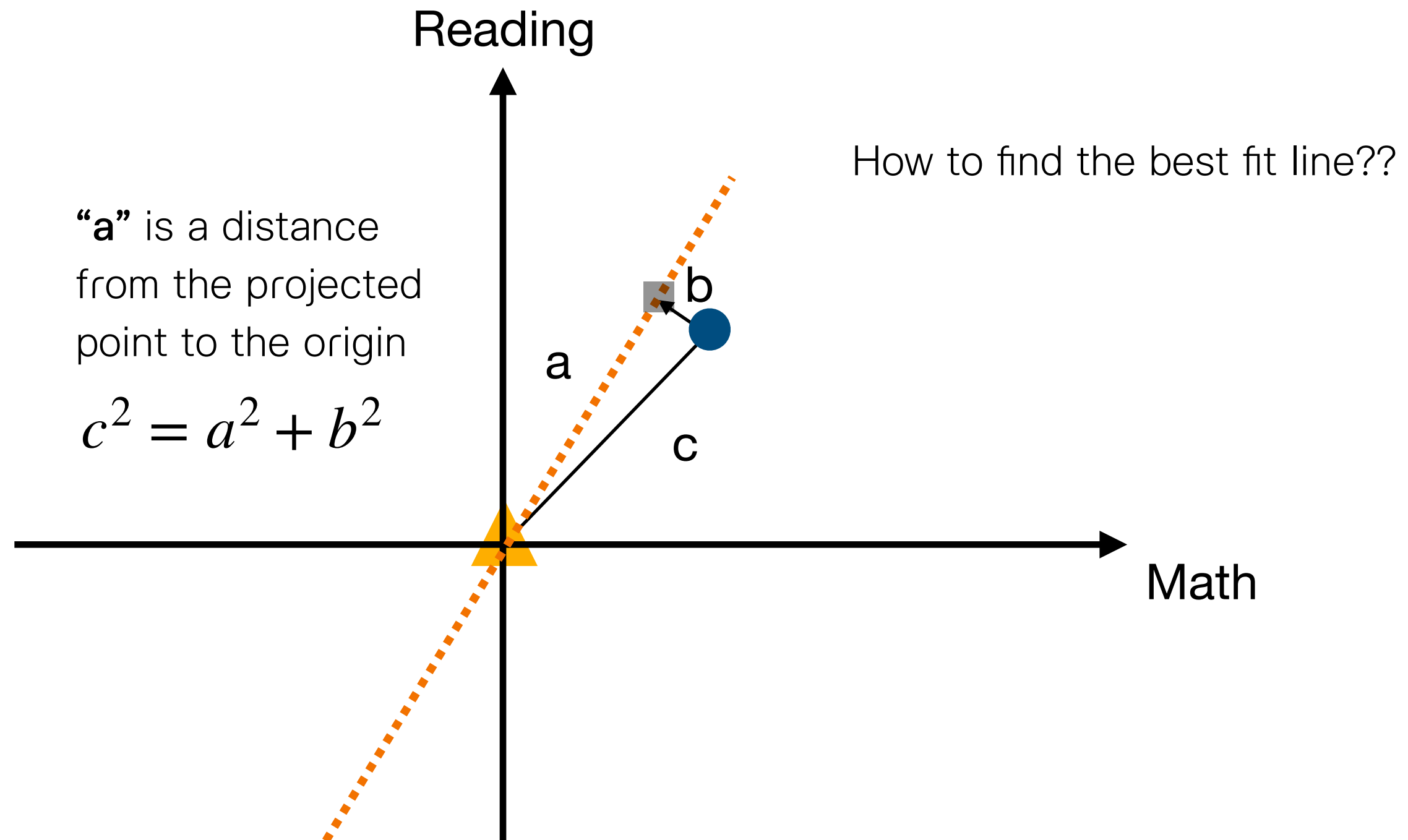
$$s^2 = \frac{\sum (y_i - \bar{y})^2}{n - 1}$$

How to find the best fit line??



PCA - concept

Fit a line to the data...

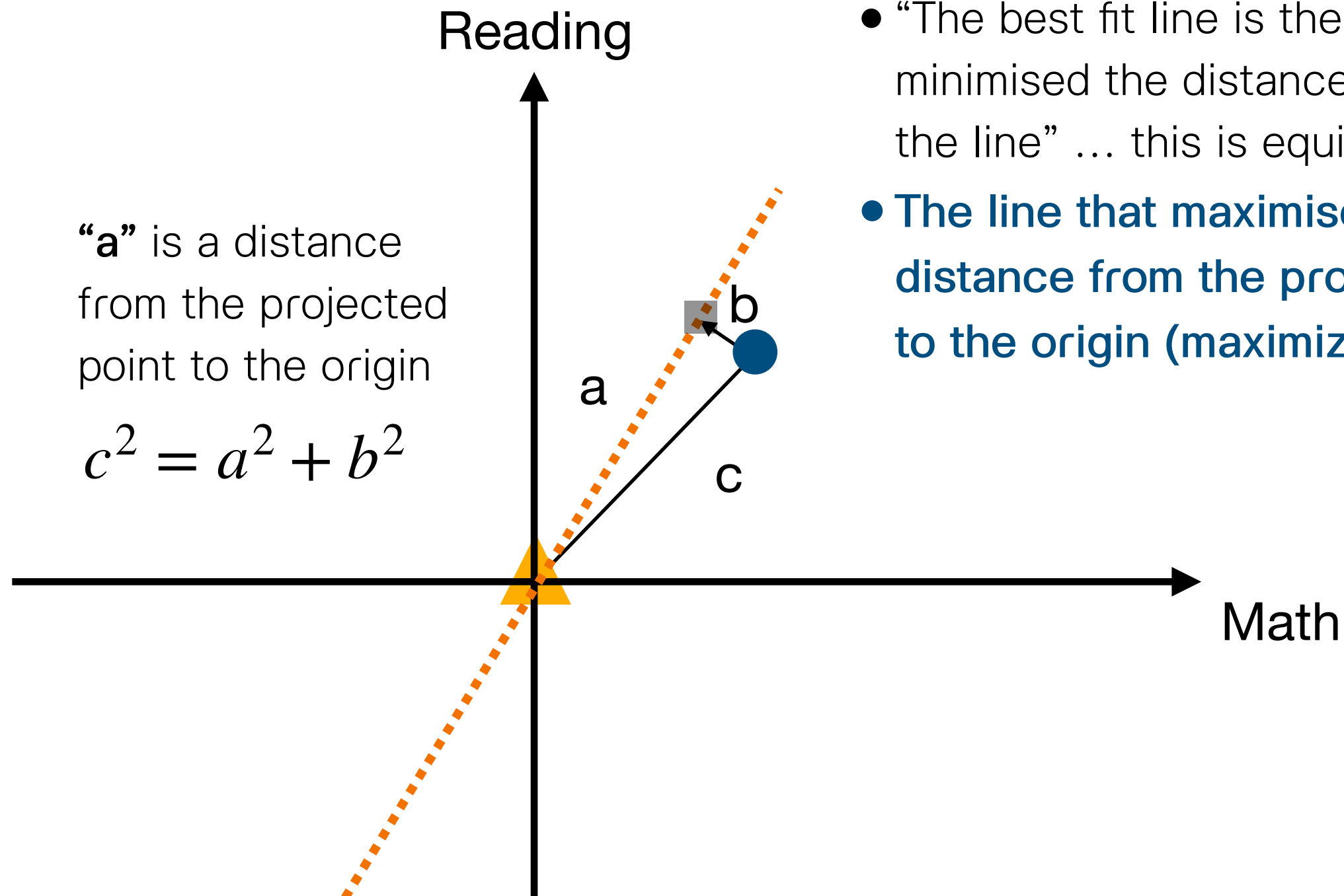


PCA - concept

Fit a line to the data...

How to find the best fit line??

- “The best fit line is the line that minimised the distances from data to the line” ... this is equivalent to....
- **The line that maximised the distance from the projected points to the origin (maximizes “a”)**

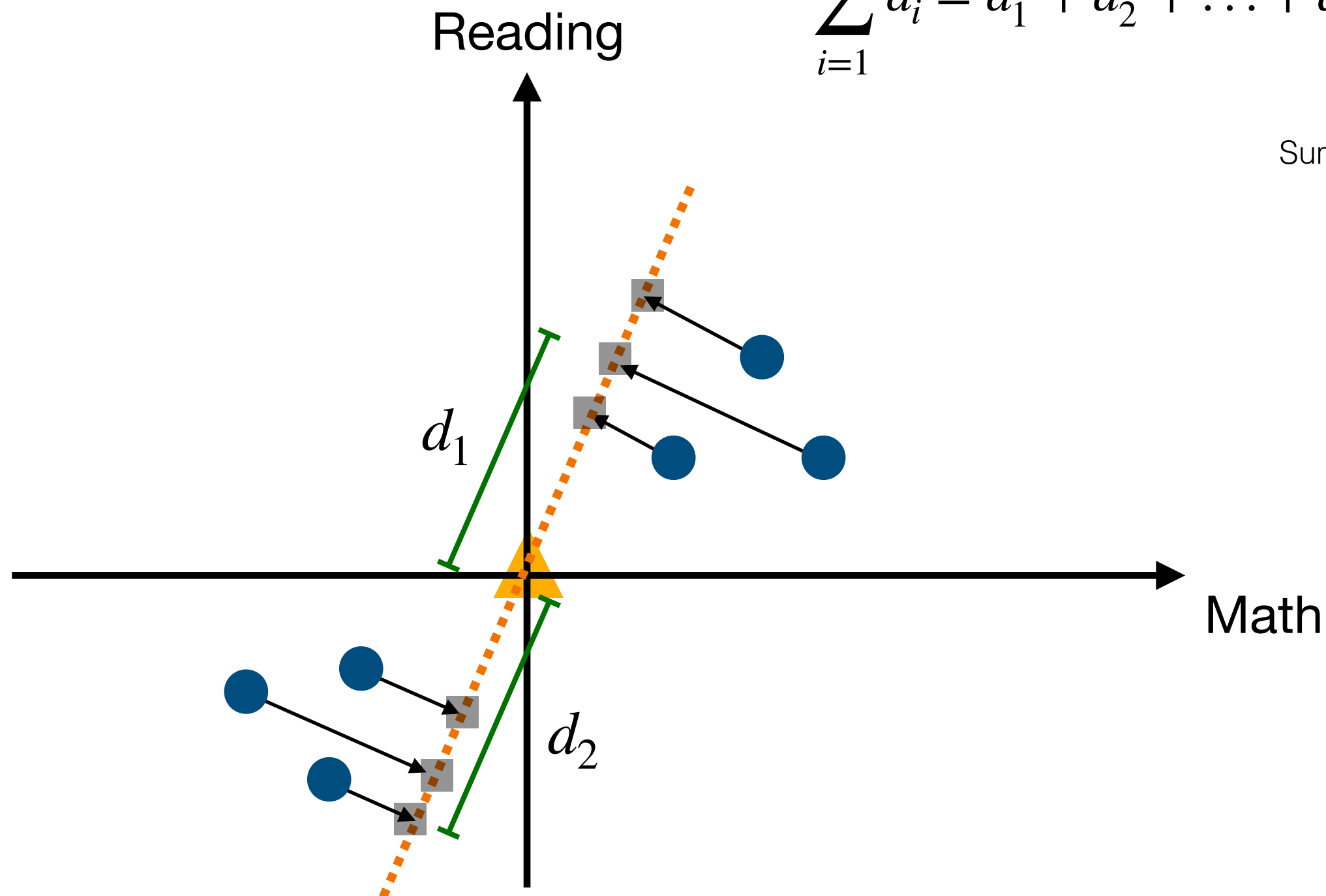


PCA - concept

Let d_i be a distance from the projected point i to the origin.

$$\sum_{i=1}^6 d_i^2 = d_1^2 + d_2^2 + \dots + d_6^2 = SS_1$$

Sum squares distance



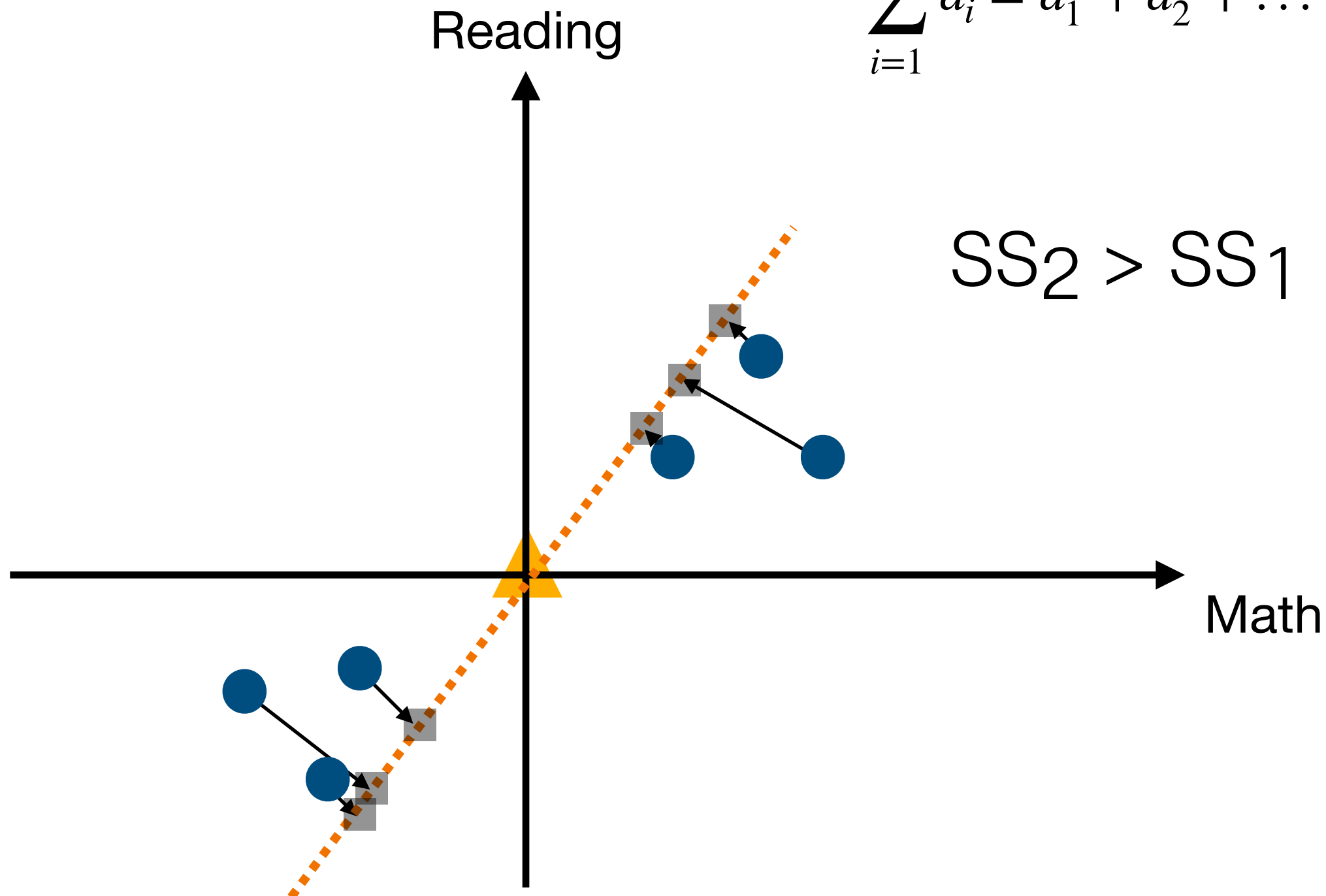
PCA - concept

Fit a line to the data...

Let d_i be a distance from the projected point i to the origin.

$$\sum_{i=1}^6 d_i^2 = d_1^2 + d_2^2 + \dots + d_6^2 = SS_2$$

$$SS_2 > SS_1$$

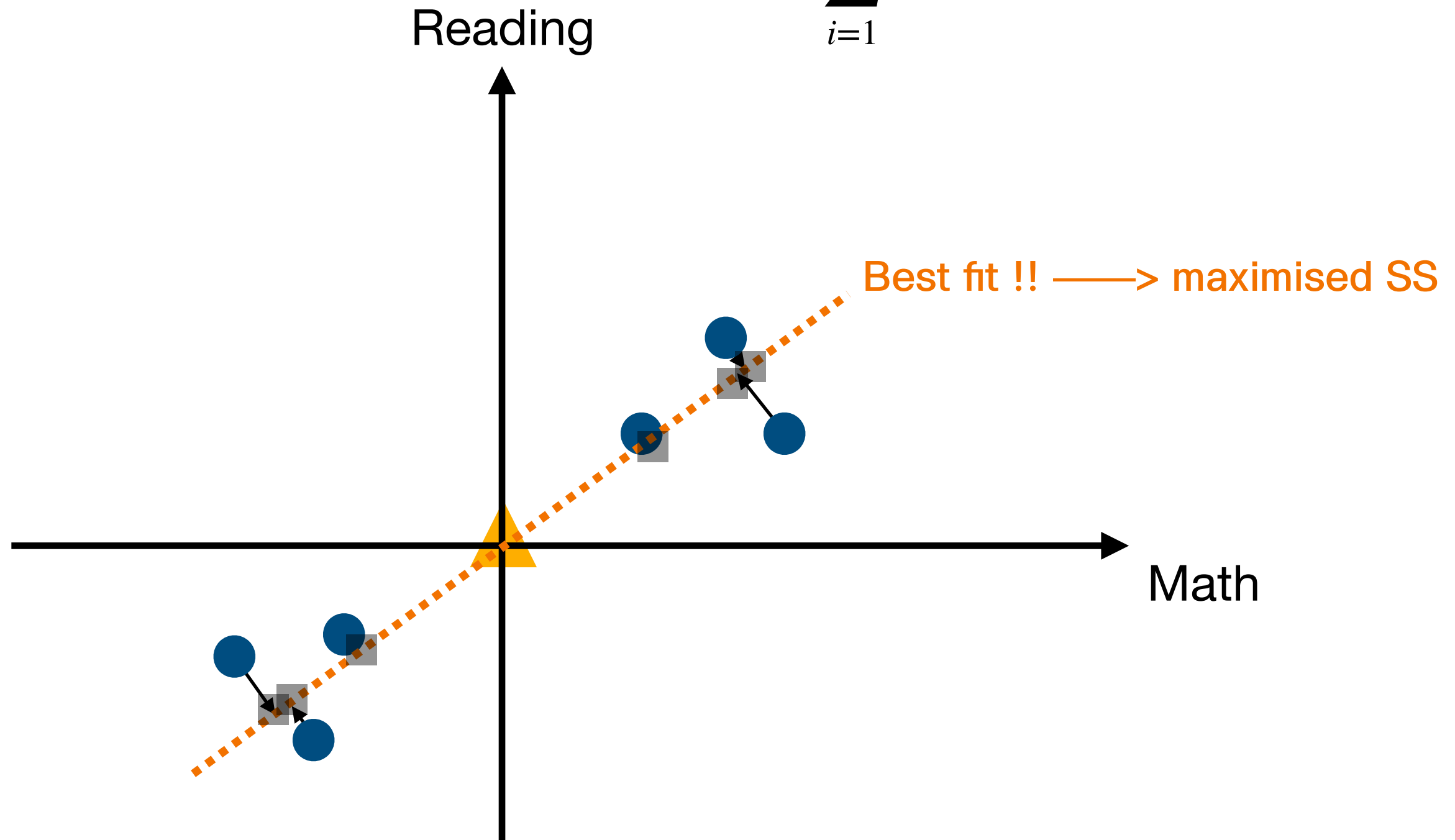


PCA - concept

Fit a line to the data...

Let d_i be a distance from the projected point i to the origin.

$$\sum_{i=1}^6 d_i^2 = d_1^2 + d_2^2 + \dots + d_6^2 = SS$$

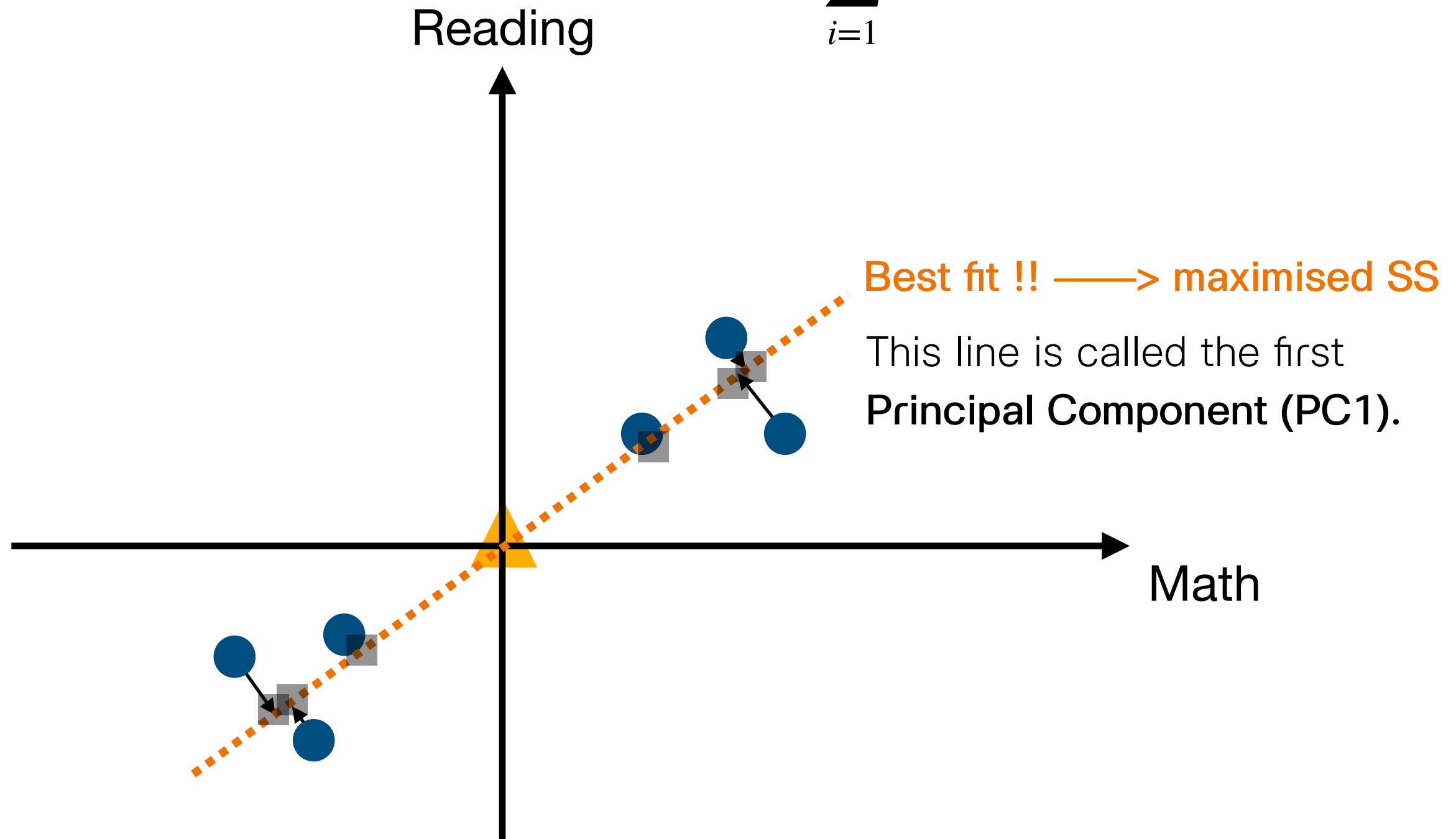


PCA - concept

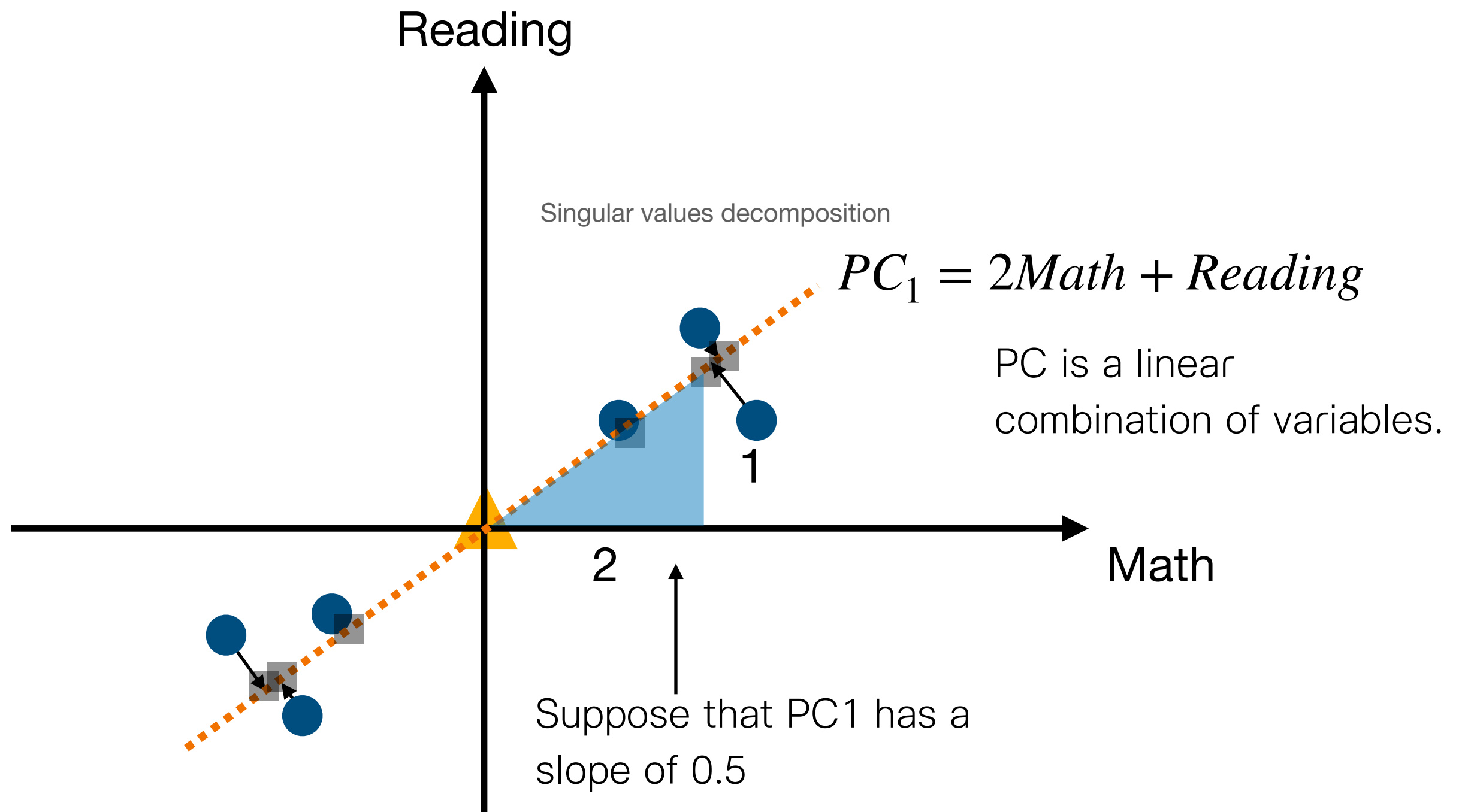
Fit a line to the data...

Let d_i be a distance from the projected point i to the origin.

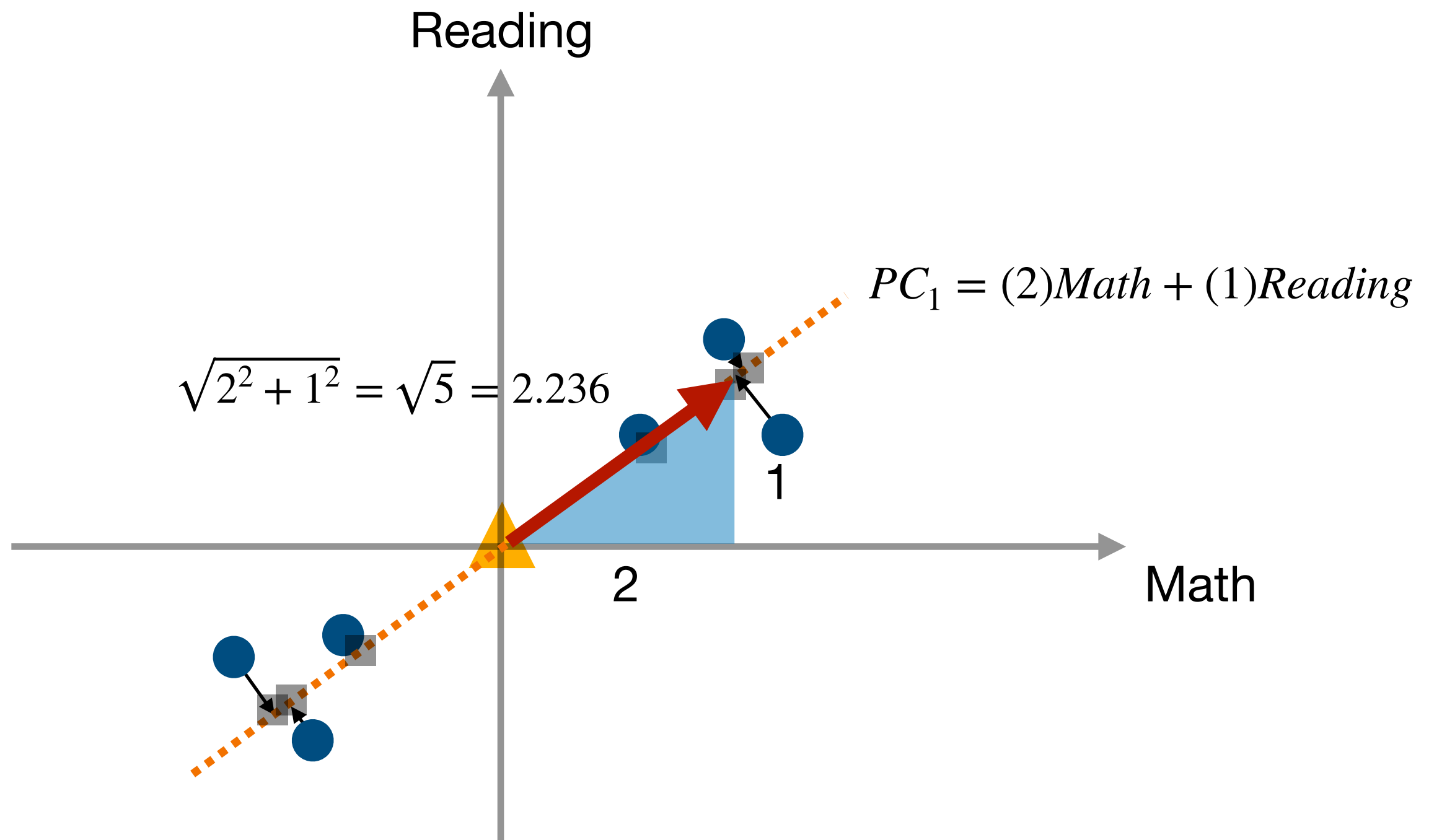
$$\sum_{i=1}^6 d_i^2 = d_1^2 + d_2^2 + \dots + d_6^2 = SS$$



PCA - concept



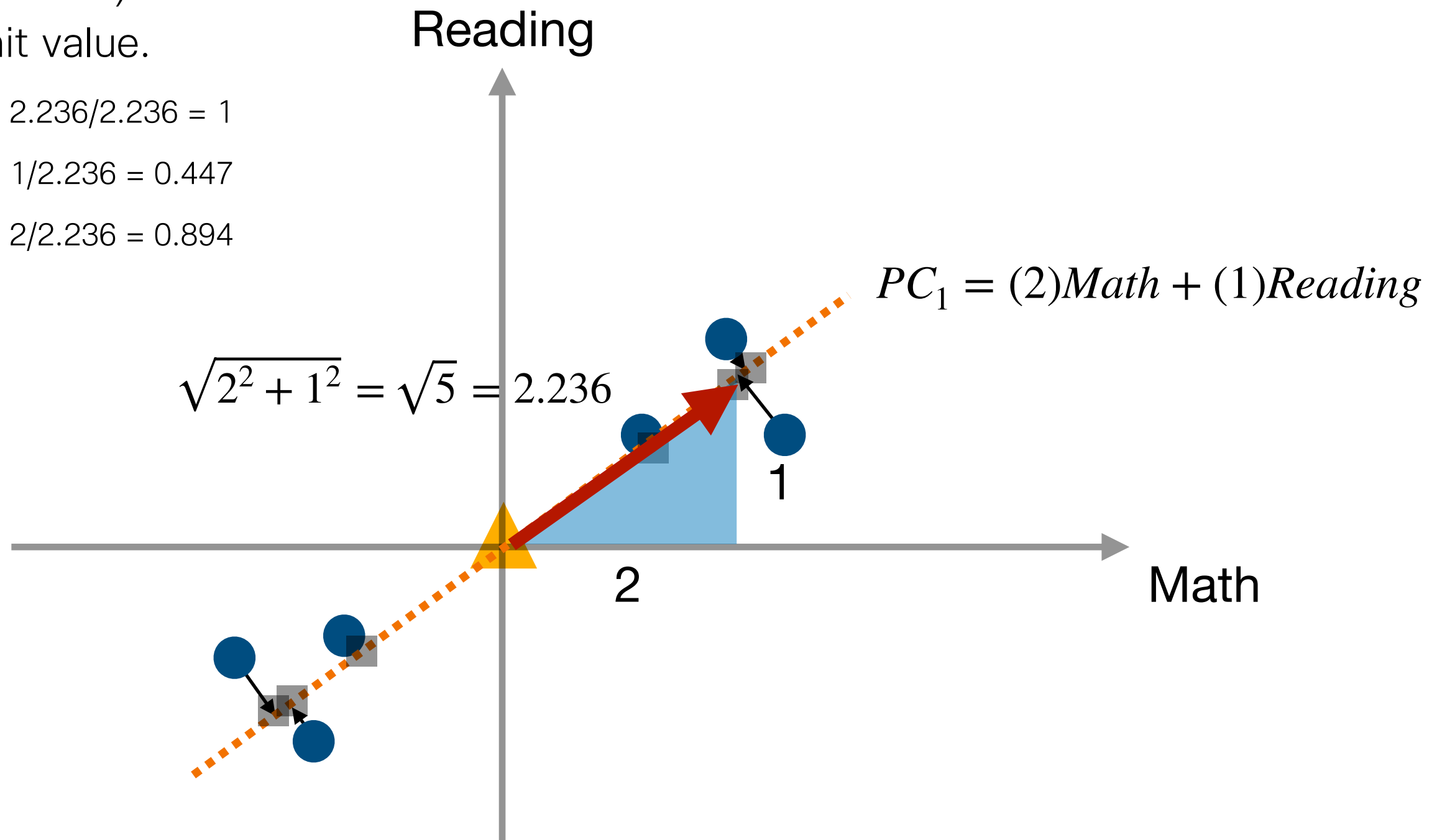
PCA - concept



PCA - concept

In PCA, we scaled the
(red line) distance into
unit value.

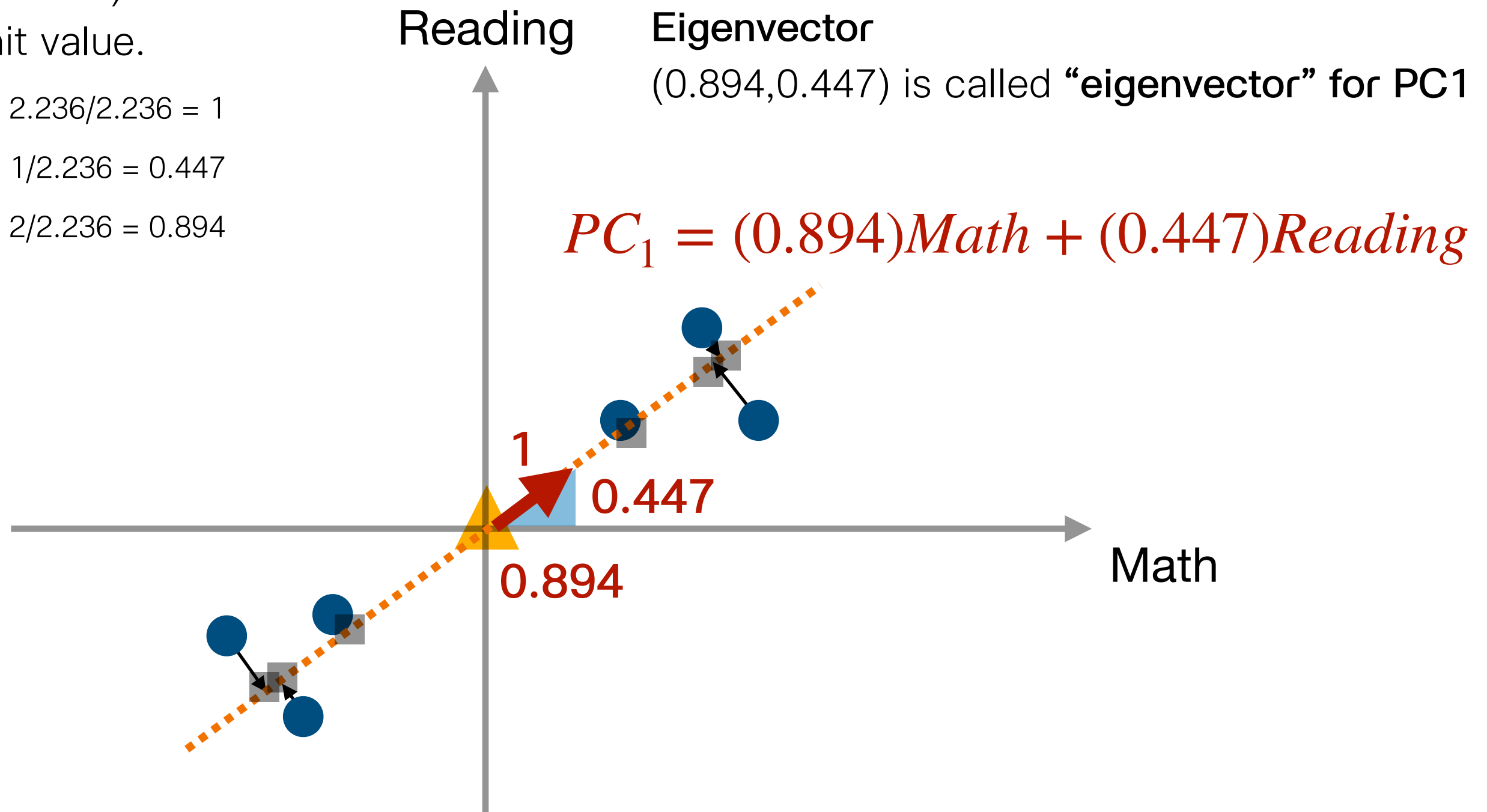
- $2.236/2.236 = 1$
- $1/2.236 = 0.447$
- $2/2.236 = 0.894$



PCA - concept

In PCA, we scaled the (red line) distance into unit value.

- $2.236/2.236 = 1$
- $1/2.236 = 0.447$
- $2/2.236 = 0.894$

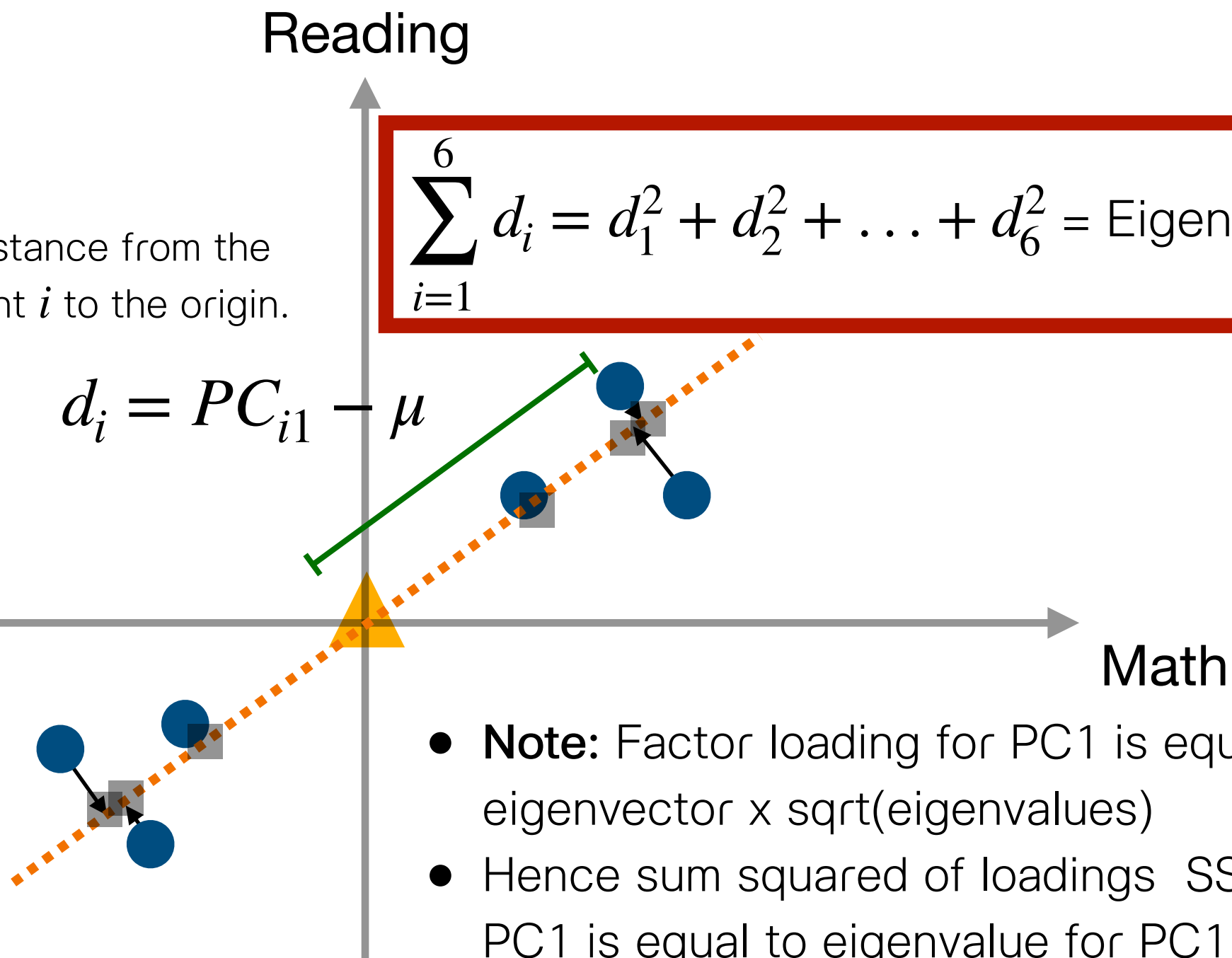


PCA - concept

Let d_i be a distance from the projected point i to the origin.

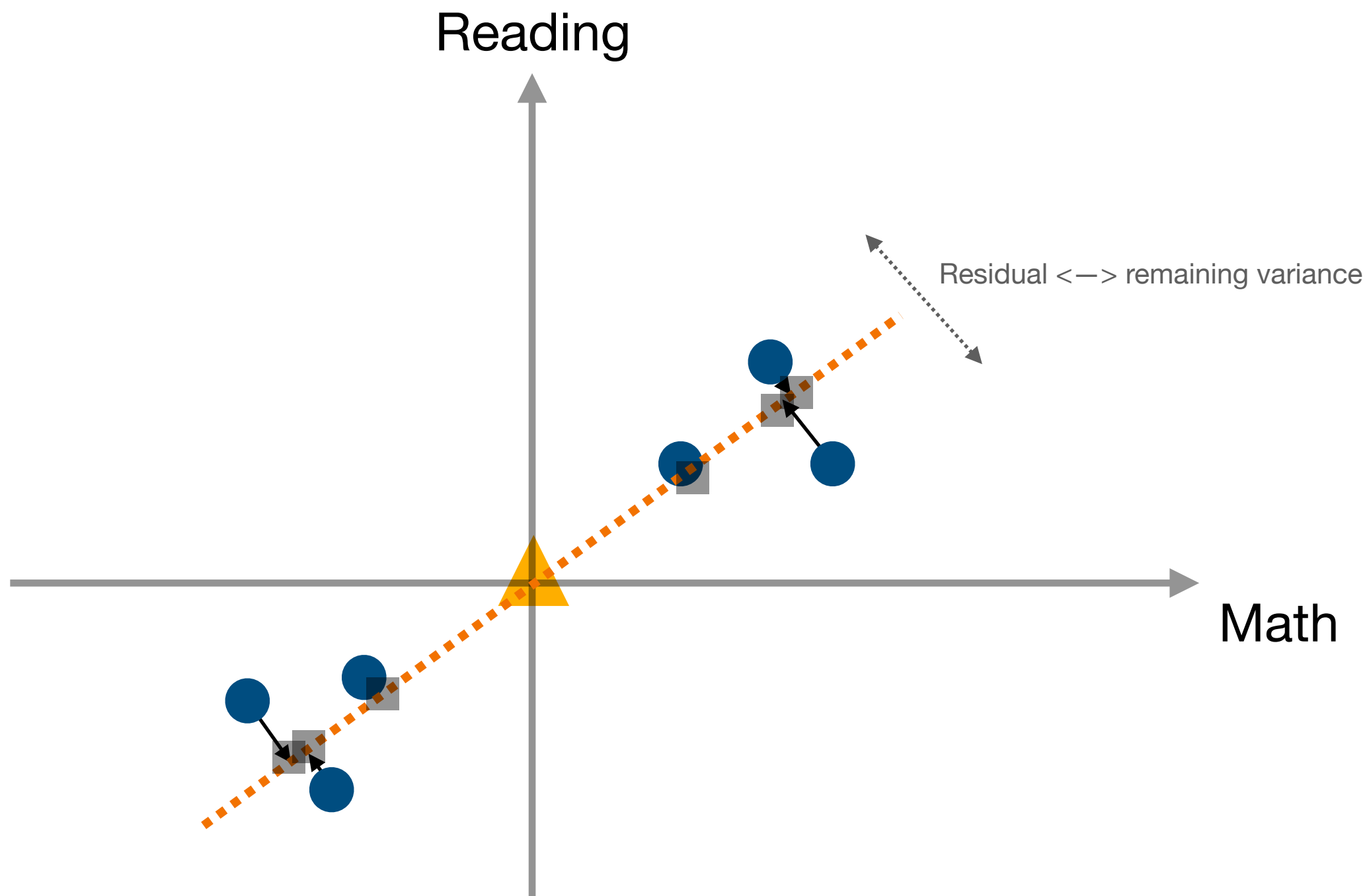
$$d_i = PC_{i1} - \mu$$

$$\sum_{i=1}^6 d_i^2 = d_1^2 + d_2^2 + \dots + d_6^2 = \text{Eigenvalue for PC1}$$

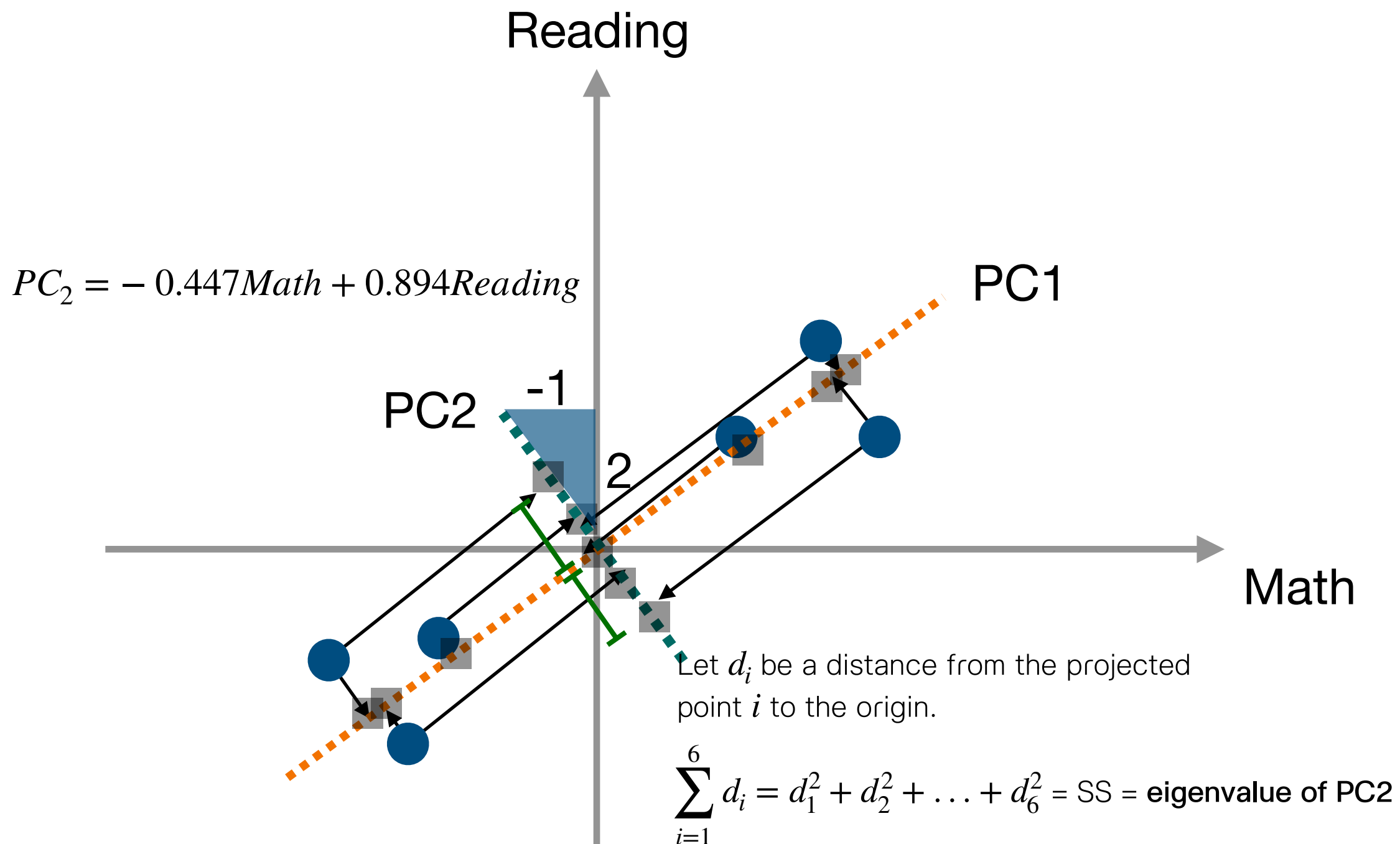


- **Note:** Factor loading for PC1 is equal to eigenvector $\times \sqrt{\text{eigenvalue}}$
- Hence sum squared of loadings SS-loading for PC1 is equal to eigenvalue for PC1

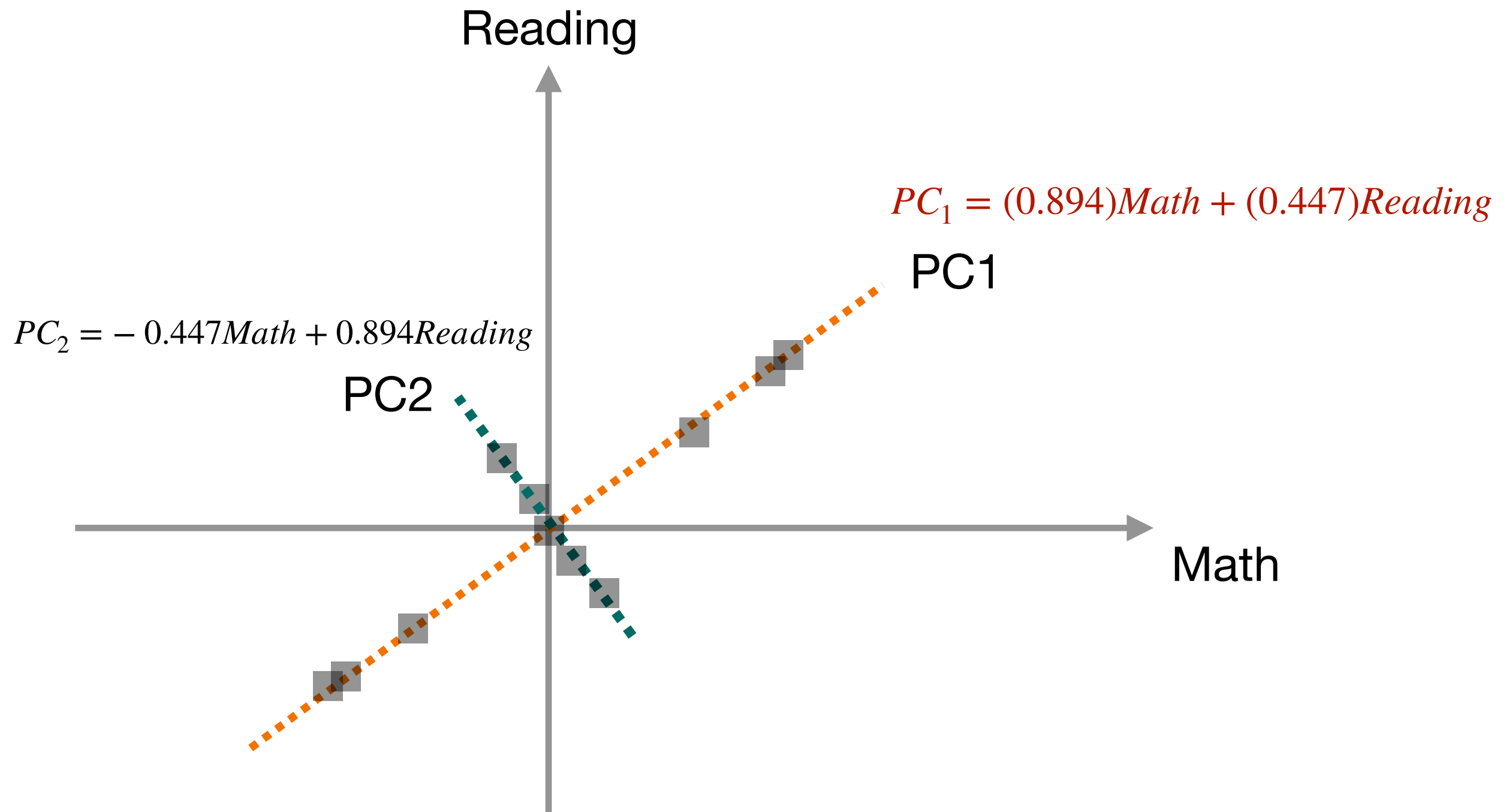
PCA - concept



PCA - concept



PCA - concept



PCA - concept

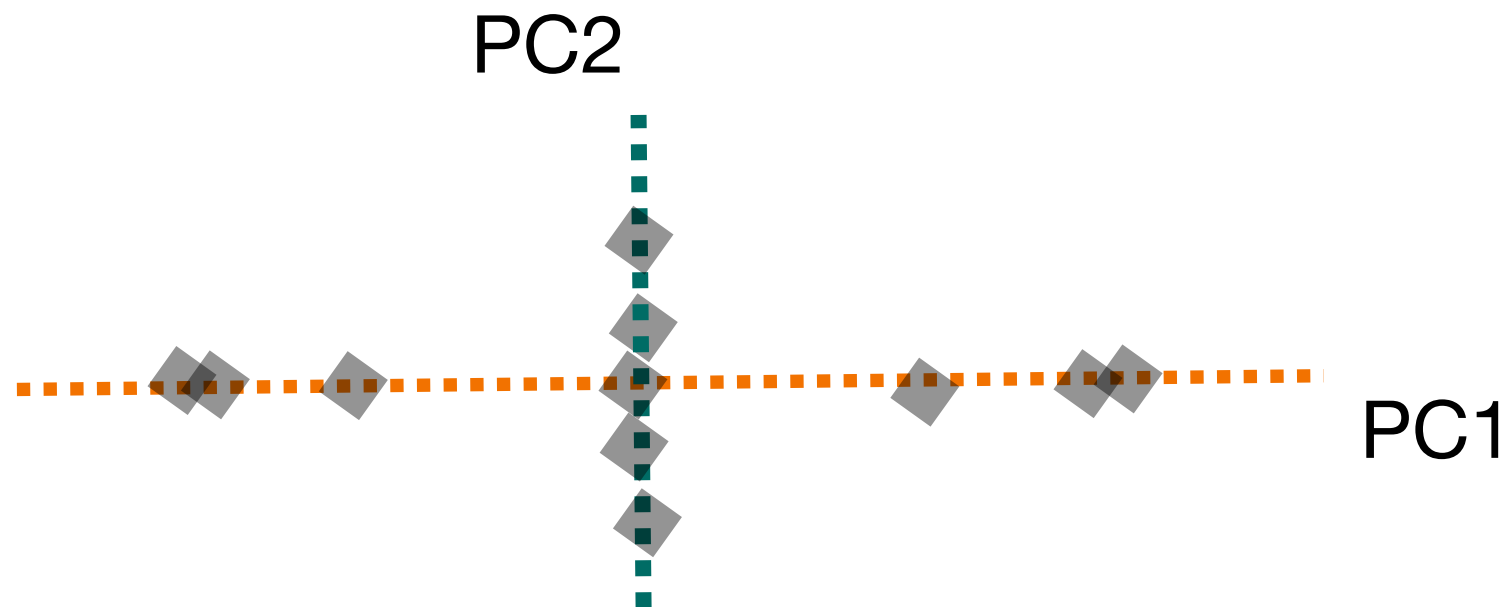


Figure not to scale.

