

★ ANOVA - Analysis Of Variance

① Cat vs. Numerical
(ex - Gender vs. Income) $\rightarrow n > 30$ - z-test
 \rightarrow T-test

② cat. vs. cat
(ex - Gender vs. Product) $\rightarrow \text{chi}^2$

③ cat vs. cat vs. cat \rightarrow ANOVA
(ex - Gender vs. Product vs. City)

When we have more than 2 columns, we use ANOVA

★ ANOVA (Deep dive)

→ Basketball players

USA - 6.6 ± 2.5

India - 6.2 ± 2.5

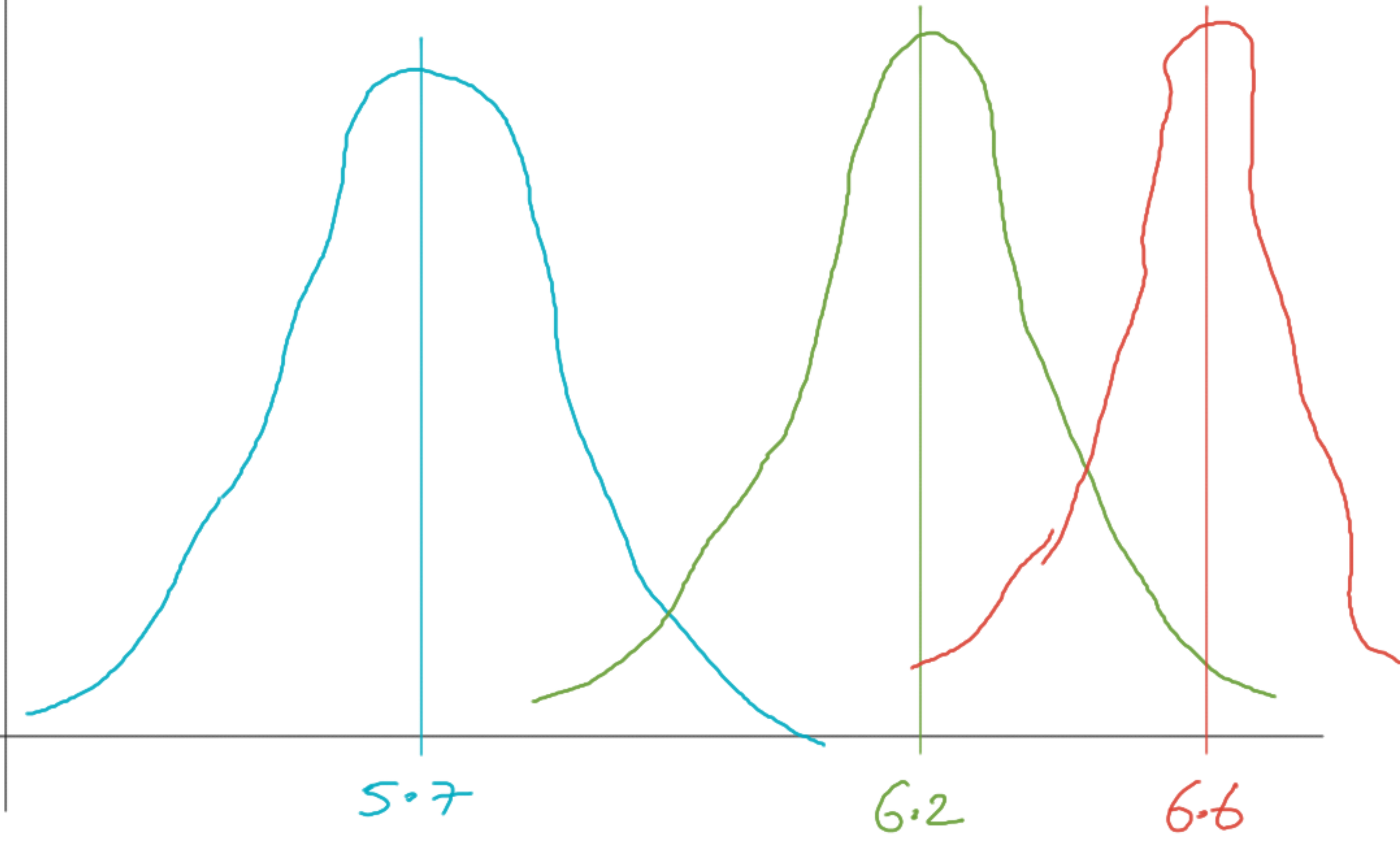
Indonesia - 5.7 ± 2.5

Question: Are mean height of Basketball players associated with groups (of country)?

Case-1

H₀: Height is groups are not associated.

H_a: They are associated.



→ If we mix-up all of them and then randomly make 3 groups - g_1, g_2 & g_3

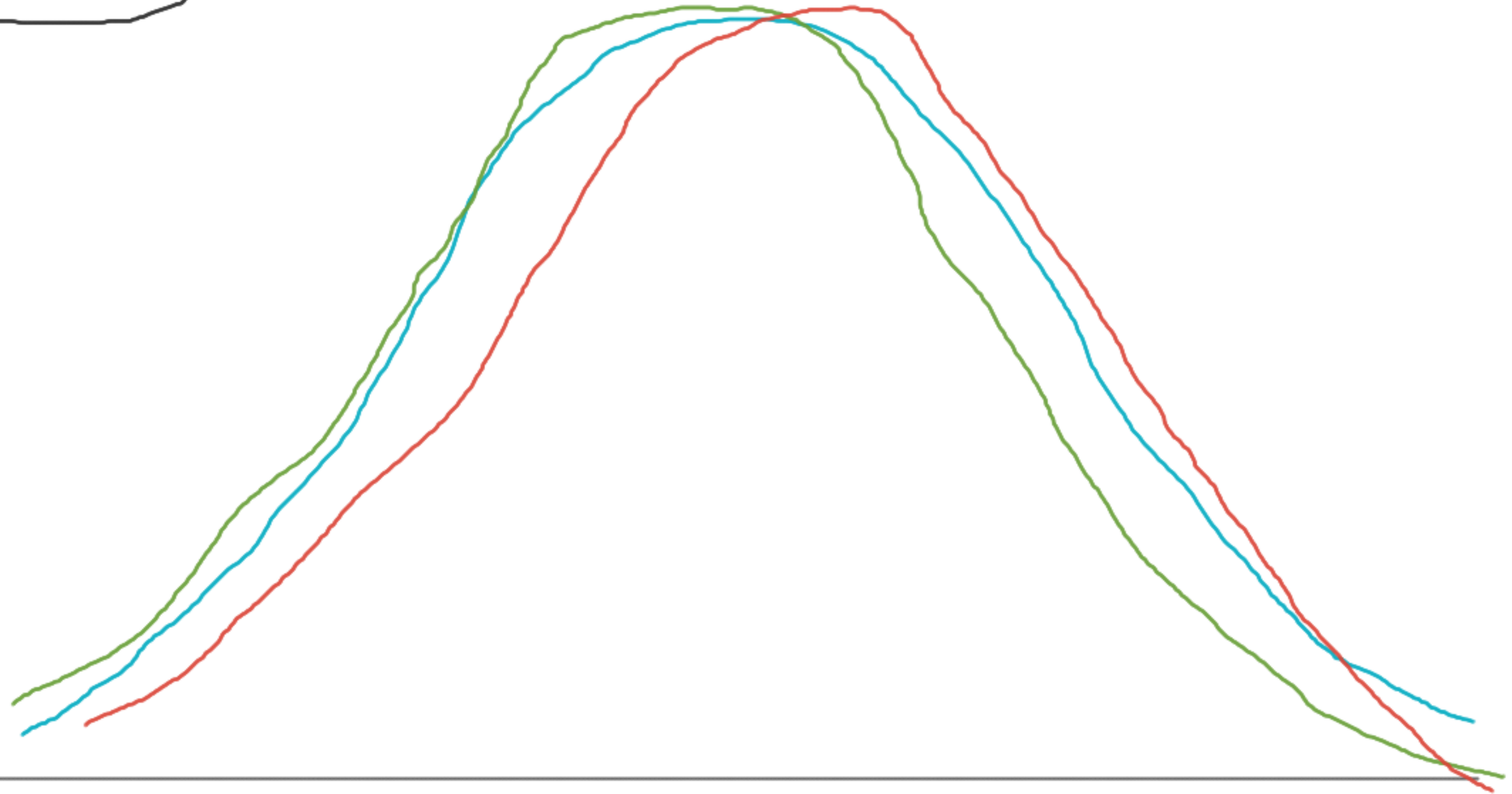
$\text{mean}(g_1)$

$\approx \text{mean}(g_2)$

$\approx \text{mean}(g_3)$

Standard deviation ↑

Case-2



$$F = \frac{\text{Variance between the groups}}{\text{variance within the groups}}$$

Case-1: variance between the groups is high &
variance within the groups is low. \therefore F-ratio is high
 \therefore p-value is lower \therefore Reject H_0

Case-2: is reverse. \therefore F-ratio is lower hence p-value is
high \therefore Fail to reject H_0

★ Assumptions of ANOVA

① Data must be Gaussian

↳ Plot the histogram

↳ Test-1: Q-Q Plot

↳ Test-2: Shapiro

② Rows must be independent of each other

③ Variance within the groups must be equal.

If these assumptions are not followed then we use Kruskal Test

Logic behind Q-Q plot :