

# exploratory data analysis and models on the epi dataset

date: 2025-10-13

## dataset and choices

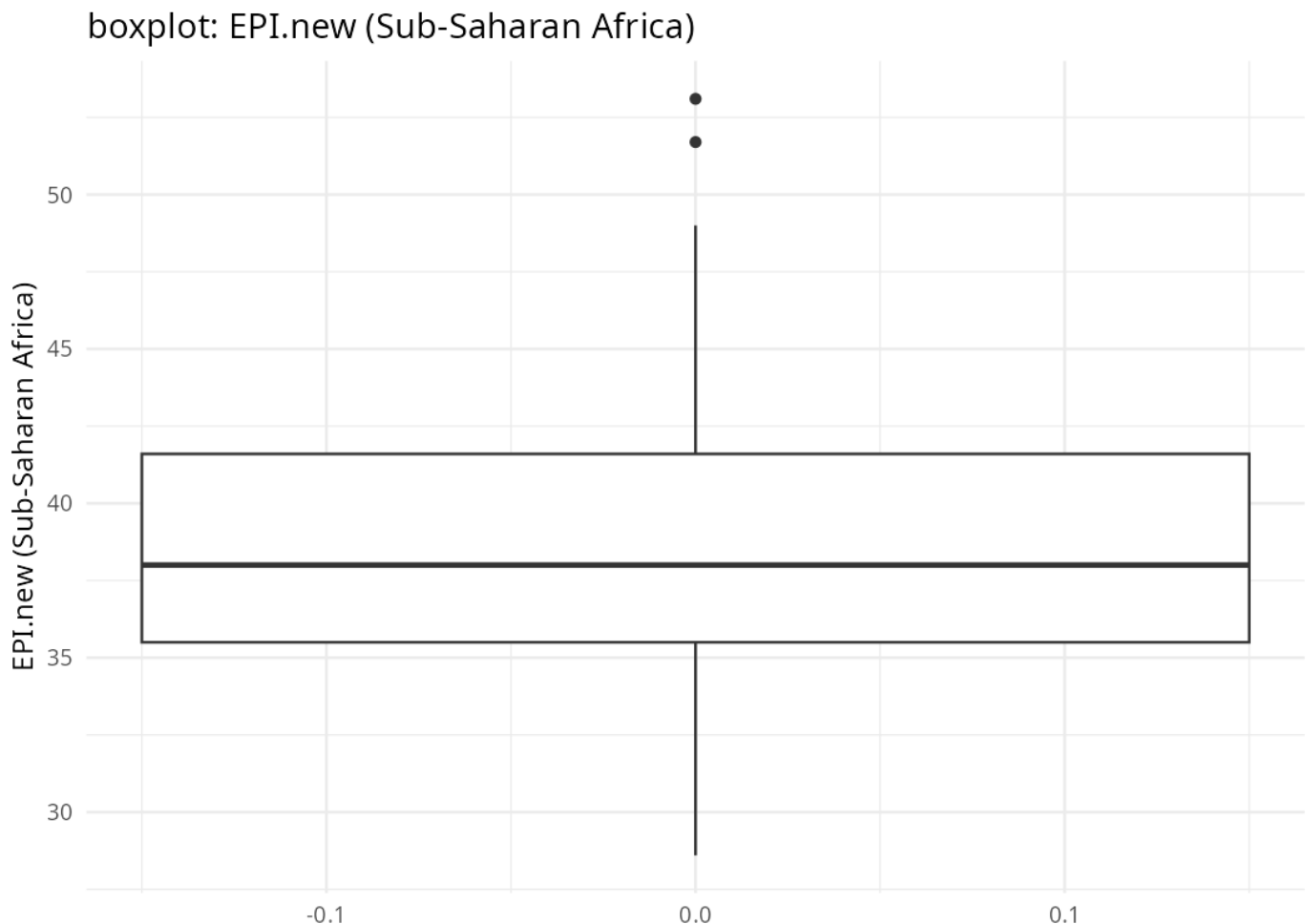
---

- **file:** epi\_results\_2024\_pop\_gdp\_v2.csv
- **region column:** region
- **response var:** EPI.new
- **regions:** Sub-Saharan Africa vs Latin America & Caribbean

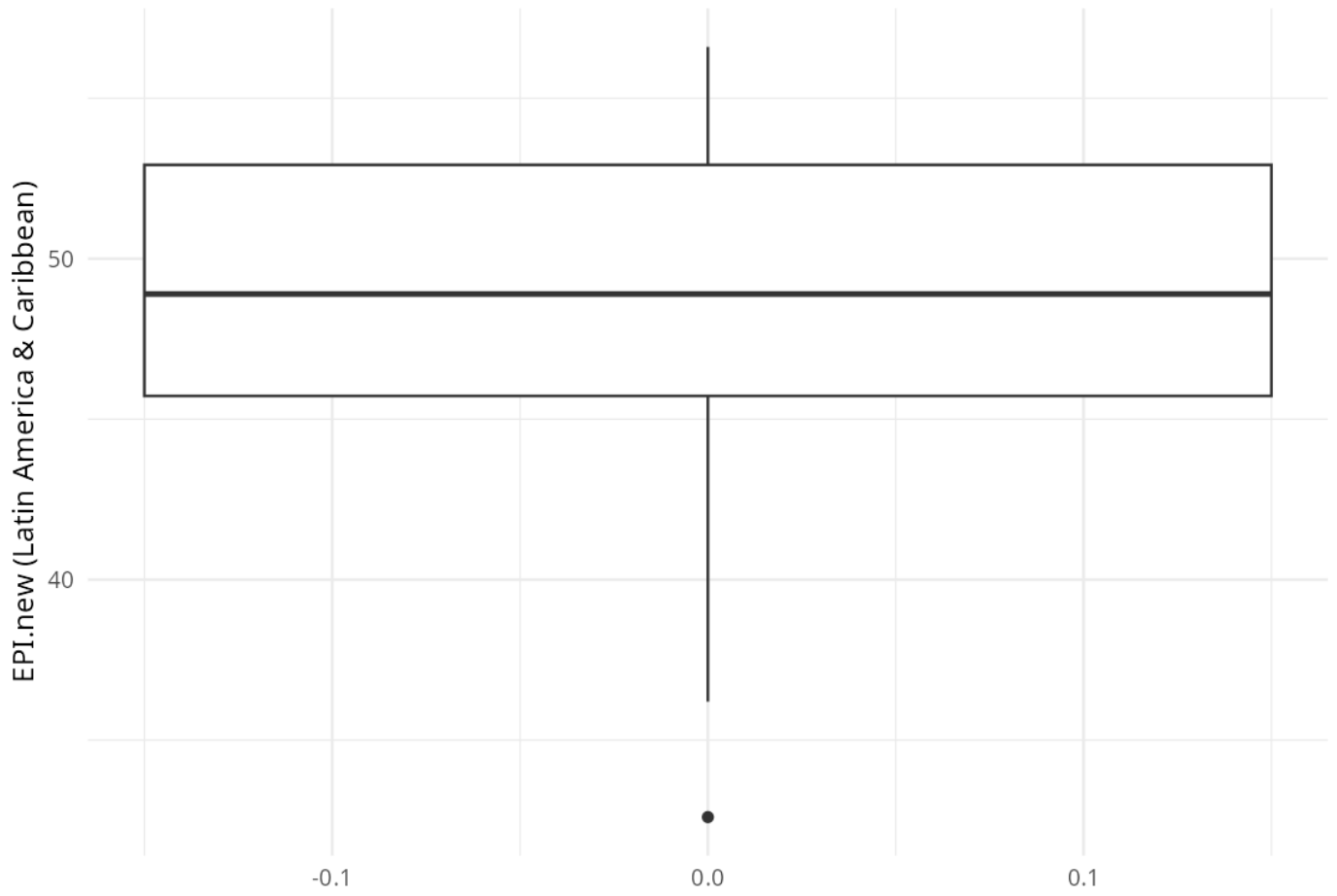
## 1) variable distributions

---

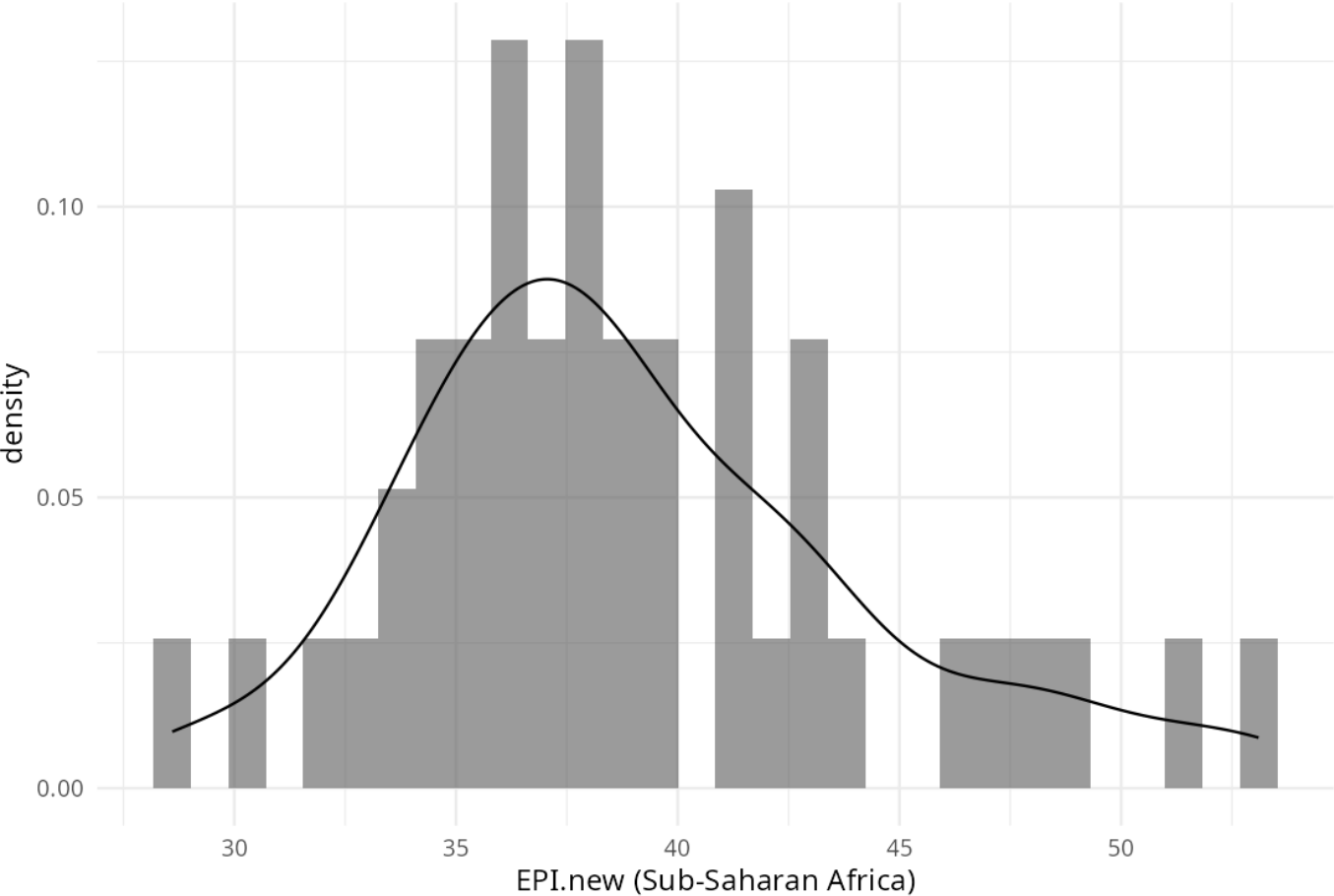
### 1.1 boxplots and histograms (with density!)



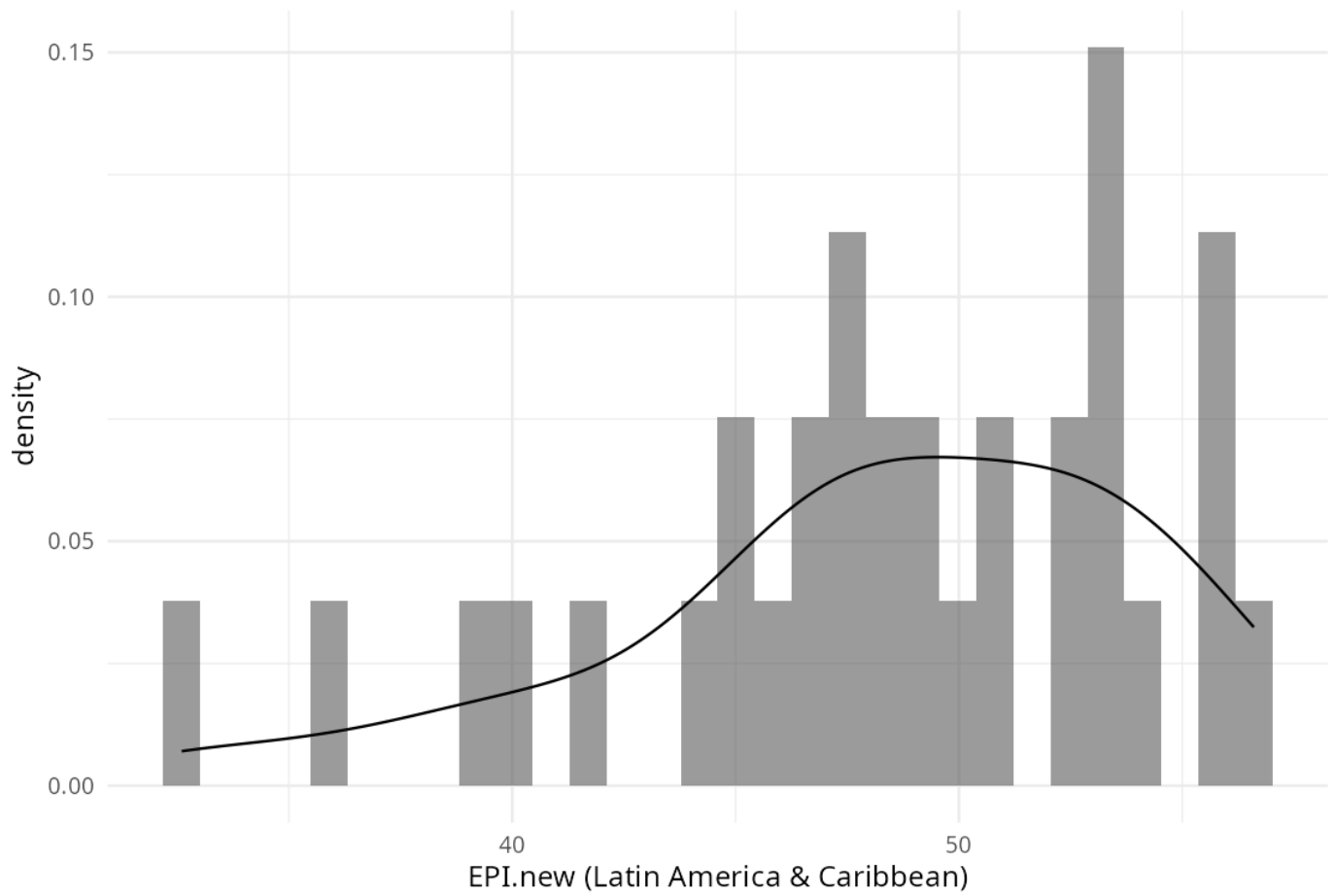
boxplot: EPI.new (Latin America & Caribbean)



histogram + density: EPI.new (Sub-Saharan Africa)

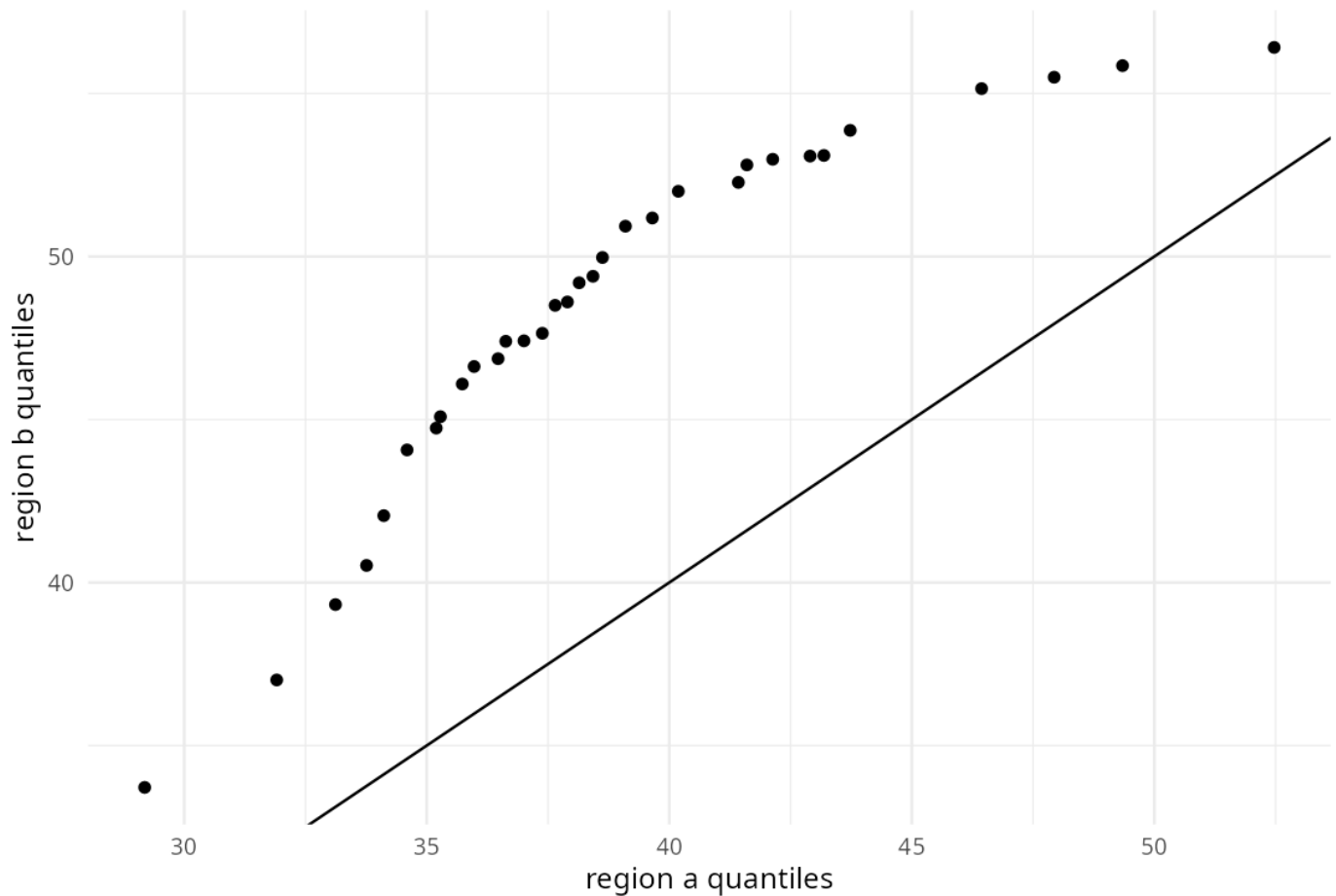


histogram + density: EPI.new (Latin America & Caribbean)



## 1.2 qq plot (two-sample)

qq plot: Sub-Saharan Africa vs Latin America & Caribbean



## 2) linear models

---

full:  $EPI.new \sim gdp$

full:  $EPI.new \sim gdp + population$

### 2.2 same models on one region (comparison)

on region Sub-Saharan Africa, the better model is **region Sub-Saharan Africa:  $EPI.new \sim gdp + population$**  ( $r^2=0.361$ ,  $aic=265.4$ ,  $bic=272.7$ ).

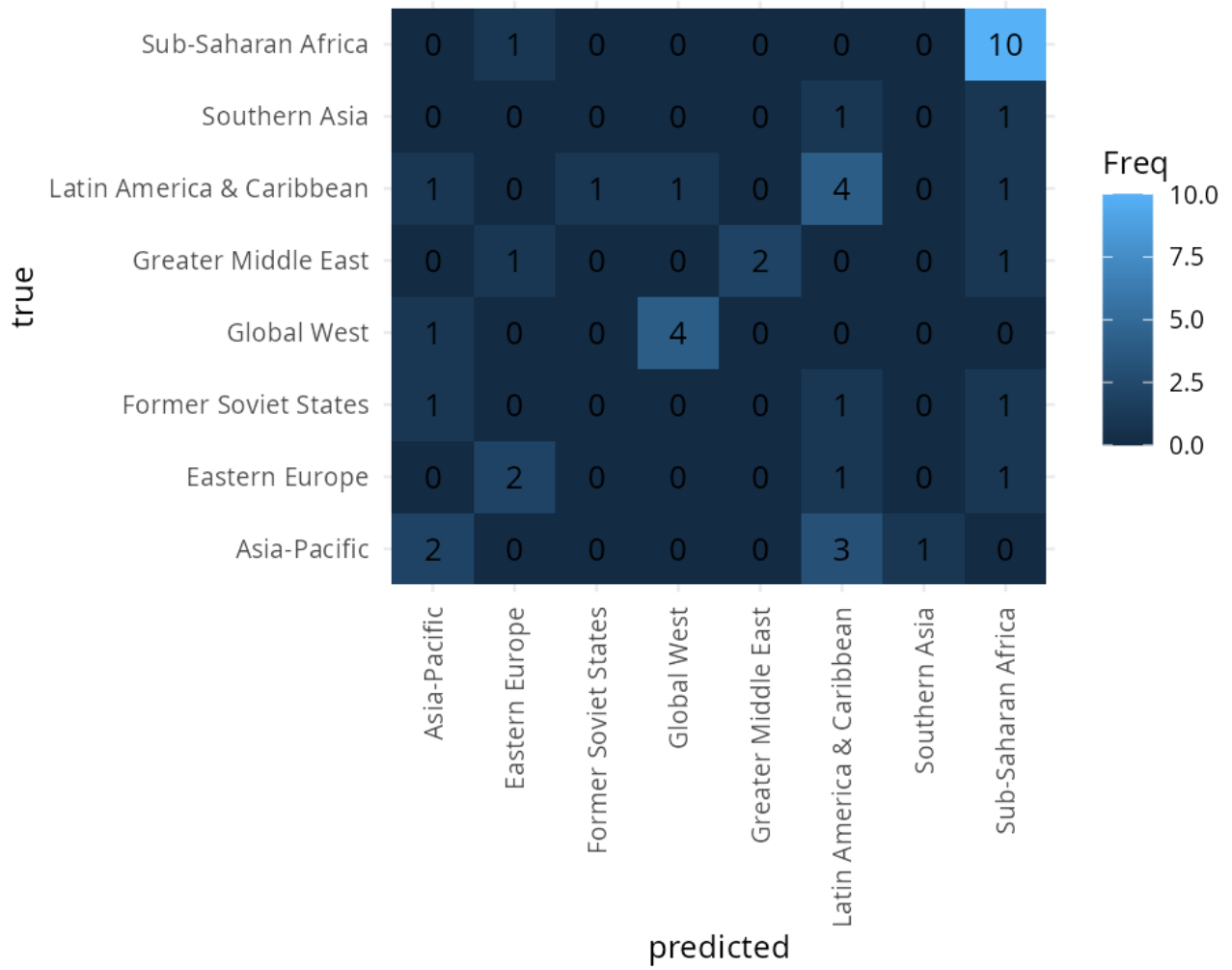
## 3) classification (knn, label = region)

---

### model A

- **k: 5** | **accuracy: 0.5581** | **test n: 43** variables: `c("AGR.new", "AIR.new", "APO.new")`

confusion matrix (k=5) – model A



**model B**

- **k:** 5 | **accuracy:** 0.5116 | **test n:** 43 variables: c("BCA.new", "BDH.new", "CBP.new")

confusion matrix (k=5) - model B

