

## Project 2: High-Dimensional Linear Models and Convergence in Economic Growth

Due: Sunday November 2<sup>nd</sup>, 2025, at 22:00

### 1 Economic Growth and Convergence

A central question in growth theory is whether economies exhibit some form of convergence. Specifically, do initially poor countries tend to grow faster than initially rich countries, thus ending up on a common path? Barro [1991] explores this question by regressing the average annual growth rate of GDP per capita ( $g_i$ ) in country  $i$  on the log of initial GDP per capita ( $y_{i0}$ ) and a vector of control variables ( $\mathbf{z}_i$ ), having in mind a relationship of the form

$$g_i = \beta y_{i0} + \mathbf{z}_i \boldsymbol{\gamma} + u_i, \quad (1)$$

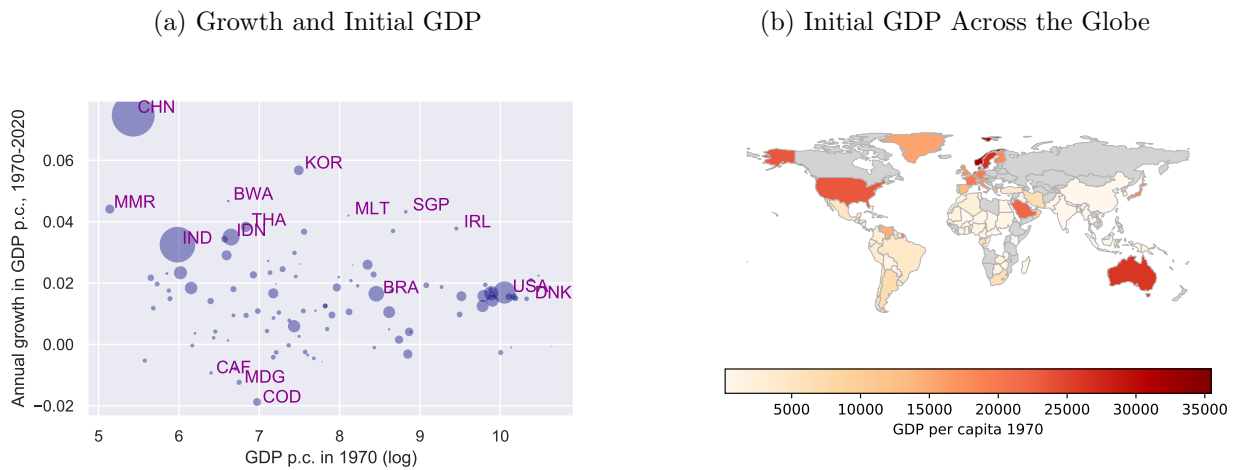
with  $u_i$  being an econometric error term (unobservable). Figure 1a shows a scatterplot of the data with  $g_i$  on the vertical axis and  $y_{i0}$  on the horizontal axis, which indicates a seemingly clear negative unconditional correlation. That is, it seems that the raw associations in the data are consistent with a negative value of  $\beta$ , which is referred to as “beta convergence” or “catch-up growth.”

However, Figure 1b shows initial GDP on a world map, from which one may be inspired to start thinking of potential omitted variables, that might simultaneously explain initial GDP and subsequent growth. For instance, although many African and Asian countries started from approximately the same levels, the subsequent growth trajectories have been remarkably different.

A key challenge in estimating  $\beta$  is therefore what to include in  $\mathbf{z}_i$  in order to avoid omitted variable bias. Over time, researchers have emphasized a long list of possible controls, including geographic factors (like temperature, disease, agricultural suitability, and natural resources), but also historical determinants like colonialization, genetic diversity, culture, and economic institutions (broadly speaking). For example, a country may have been poor in 1970 due to a hazardous climate, which then also explains why that country remains poor and does not catch up during 1970–2020.

With as many (or more) candidate regressors as there are countries, the empirical researcher faces a high-dimensional estimation problem when attempting to estimate  $\beta$ : Ideally, we want to avoid restricting the set of possible control variables ex ante, and instead let the data tell us which controls are most important for the purpose of explaining economic

Figure 1: Growth and Convergence



*Note:* In panel a, each dot is a country, and the size of the dot is proportional to the initial population. Names are provided for a few selected countries. In panel b, countries are colored by initial GDP per capita (1970). Grey colors indicate missing values.

growth. However, we also recognize the selection problem: With so many potential controls, we are forced to choose some and leave out others. (Even high-dimensional methods break down eventually.)

## 2 Data

The dataset, `growth.csv`, contains data on 214 countries for a long list of variables (see Appendix A). For the intended outcome variable,  $g_i$ , you should use the variable `gdp_growth`, and the key explanatory variable of interest,  $y_{i0}$ , is `lgdp_initial`. As the appendix table shows, `gdp_growth` is only available for 102 countries, so this selection implies an upper bound on the number of observations.

You must choose which variables to include in your analysis in addition to these two. For instance, it is common to include the investment rate, `investment_rate`, as a control, but what additional controls you permit and how to include them is up to you. Aim to strike a balance between permitting as many as possible and avoiding variables that restrict your sample size too heavily (although such a restriction could in itself be telling a story).

## 3 Assignment

Test whether the data is consistent with the theory of beta convergence. In doing so, you must treat  $\mathbf{z}_i$  as high dimensional.

## 4 Hints

- (1) You may start by discarding some of the available variables and never considering those in your project. For example, you could write your entire paper only considering the geographical controls. You should then proceed by treating that list of variables in a high-dimensional paradigm. If you discard something, provide your reasoning.
- (2) When using an estimation procedure, carefully discuss the assumptions required to derive the estimator and establish properties thereof. Assess whether these assumptions are likely to be satisfied in the current empirical setting. (Don't just copy the math; relate to the current setting.) If not, what are the consequences for the estimator in question (and your results)? Strive to provide a real-world example of behavior that might invalidate a given assumption, carefully linking the behavior or mechanisms to the mathematical symbols in the model.
- (3) If you rationalize several model specifications and associated estimates, discuss which one seems the most appropriate and justify your decision (e.g. based on formal testing).
- (4) Be precise about the statistical tests you use for testing various hypotheses. Explain which null hypothesis you are testing and the alternative you are testing against, how the test statistic is constructed, the decision rule you employ, and the conclusion you reach. If a variance (matrix) has been estimated, discuss the assumptions invoked for consistency. If several choices are possible, justify your choice.

## 5 Formal Requirements

- You must hand in a report that presents the econometric model, presents your estimation results and results of formal statistical tests (including interpretation and statements on economic and statistical significance), and discusses the potential weaknesses of the model, data and approach. If you present many estimates of the same parameters (e.g., estimators based on different assumptions, or varying the controls or sub-sample used), it may be helpful to present the estimates together in one table to facilitate comparison.
- The report must be written in English using an academic language and structure and uploaded to FeedbackFruits via Absalon as a single PDF file.
- The report must be structured as at most five pages of main text (including mathematics) followed by at most two pages of output. Use the following formatting:

- For the main text (and mathematics), you must use fontsize = 12p, line spacing = 1.5, and page margins = 2.5 cm (as used in this document). Note that this may differ from your idea of a “normal” page.
- The output can be any (relevant) tables, graphs or images as long as they are properly formatted, labelled and readable. Place the output at the end of your report, starting on a new page. Do not worry about this order being inconvenient for the reader.
- Along with your report, you must upload a compressed zip-folder with all the Python code needed to replicate your results. Make sure that your code is transparent and runs with only minor modifications (e.g., changing relevant paths). There is no character limit on the submitted Python code.
- You are allowed (and strongly encouraged) to work in groups of up to three people. List all group members on the front page of your report in alphabetical order of surnames.
- The general assessment criteria are given on the course website in Absalon, where you can also find some tips on how to write and structure your report.

## References

- Daron Acemoglu, Simon Johnson, and James A Robinson. The colonial origins of comparative development: An empirical investigation. *American economic review*, 91(5):1369–1401, 2001.
- Daron Acemoglu, Suresh Naidu, Pascual Restrepo, and James A Robinson. Democracy does cause growth. *Journal of political economy*, 127(1):47–100, 2019.
- Quamrul Ashraf and Oded Galor. The ‘out of africa’ hypothesis, human genetic diversity, and comparative economic development. *American Economic Review*, 103(1):1–46, 2013.
- Valentina A Assenova and Matthew Regele. Revisiting the effect of colonial institutions on comparative economic development. *PloS one*, 12(5):e0177100, 2017.
- Robert J Barro. Economic growth in a cross section of countries. *The quarterly journal of economics*, 106(2):407–443, 1991.

## A Variable Labels

The table below shows the descriptions and sources for all variables in the dataset in `growth.csv`. There are 208 rows – one per country – so the column “Obs.” indicates how many of the 208 countries have non-missing values for the given variable. The column “Source” indicates where the data comes from, with the abbreviations:

- WB: World Bank,
- AG: [Ashraf and Galor \[2013\]](#),
- ANRR: [Acemoglu et al. \[2019\]](#),
- AR: [Assenova and Regele \[2017\]](#). Their data is based on [Acemoglu et al. \[2001\]](#).

Variable	Description	Source	Obs.
abslat	Absolute latitude	AG	205
africa	Africa dummy	AG	208
americas	Americas dummy	AG	208
area	Total land area	AG	208
area_ar	Arable land area	AG	196
asia	Asia dummy	AG	208
capital_growth_pct_gdp_initial	Gross capital formation in 1970 (% of GDP)	WB	97
capital_growth_pct_gdp_now	Gross capital formation in 2020 (% of GDP)	WB	123
cenlong	Geodesic centroid longitude	AG	208
code	World Bank country code	AG	214
cons00a	constraint on executive in 1900	AR	91
currentinst	IFC ease of doing business index/rank, 2012 (** our measure)	AR	155
dem	Democracy measure by ANRR	ANRR	155
demBMR	Democracy measure by BMR	ANRR	154
demCGV	Democracy measure by CGV	ANRR	150
democ00a	democracy in 1900	AR	87
democ1	democracy in first year of independence	AR	87
demreg	Average democracy in the region*initial regime (leaving own country out)	ANRR	183
distc	mean distance to coast	ANRR	159
distr	mean distance to coast or river	ANRR	159
distr	mean distance to river	ANRR	159
elevavg	Mean elevation	AG	184
elevstd	Standard deviation of elevation	AG	161
europa	Europe dummy	AG	208
excolony	=1 if was colony FLOPS definiti	AR	78

gdp_growth	Annual growth in GDP per capita, 1970-2020	WB	102
gdp_initial	GDP in 1970	WB	109
gdp_now	GDP in 2020	WB	167
gdp_pc_initial	GDP per capita in 1970	WB	109
gdp_pc_now	GDP per capita in 2020	WB	167
ginv	Gross investment as a share of GDP	ANRR	104
goldm	Natural minerals: gold	AR	159
imputedmort	imputed mortality rate from logem4 measure (=exp(logem4))	AR	78
imr95	infant mortality rate (1995)	AR	60
investment_rate	Capital formation (% of GDP per year, avg. of available years 1970-2020)	WB	179
iron	Natural mineral: iron	AR	159
kgatr	Percentage of population living in tropical zones	AG	160
landlock	=1 if landlocked	AR	163
leb95	life expectancy at birth (1995)	AR	60
legor_fr	French legal origin dummy	AG	202
legor_uk	British legal origin dummy	AG	202
lgdp_initial	GDP per capita in 1970 (log)	WB	109
lh_bl	Percentage of population with tertiary education (Barro- Lee)	ANRR	143
ln_yst	Log [Neolithic transition timing]	AG	164
ln_yst_aa	Log [Neolithic transition timing (ancestry adjusted)]	AG	158
logem4	log of mortality rate (IV)	AR	87
lp_bl	Percentage of population with at most primary education (Barro-Lee)	ANRR	143
lpop_initial	Population in 1970 (log)	WB	199
ls_bl	Percentage of population with at most secondary education (Barro-Lee)	ANRR	143
lt100km	amount of territory within 100 km of the coast (McArthur and Sachs)	AR	61
malfal	Percentage of population at risk of contracting malaria	AG	164
marketref	Index of market reforms (1960)	ANRR	136
mortality	mortality measure (deaths per 1000 soldiers) from Appendix Table A2, AJR	AR	61
oceania	Oceania dummy	AG	208
oilres	oil reserves	AR	154
pcatholic	Share of Roman Catholics in the population	AG	204
pd1	Population density in 1 CE	AG	155
pd1000	Population density in 1000 CE	AG	177
pd1500	Population density in 1500 CE	AG	184
pd1500	Population density in 1500 CE	AG	184
pdiv	Predicted genetic diversity	AG	207
pdiv_aa	Predicted genetic diversity (ancestry adjusted)	AG	164

pdivhmi	Mobility index-predicted genetic diversity	AG	139
pdivhmi_aa	Mobility index-predicted genetic diversity (ancestry adjusted)	AG	132
pmuslim	Share of Muslims in the population	AG	204
polity	Polity (measure of democracy)	ANRR	122
polity2	Polity 2 (measure of democracy)	ANRR	122
pop1	Population in 1 CE	AG	155
pop1000	Population in 1000 CE	AG	177
pop1500	Population in 1500 CE	AG	184
pop_growth	Annual growth in population, 1970-2020	WB	198
population_initial	Population in 1970	WB	199
population_now	Population in 2020	WB	198
pothor	Share of other religions in the population	AG	201
pprotest	Share of Protestants in the population	AG	201
precip	Precipitation	AG	184
rough	Terrain roughness	AG	184
silv	Natural mineral: silver	AR	159
suitavg	Land suitability for agriculture	AG	155
suitgini	Land suitability Gini	AG	160
temp	Temperature	AG	184
tropicar	% land area in geographical tropics	ANRR	159
uvdamage	Ultraviolet exposure	AG	207
yellow	=1 if vector yellow fever present today	AR	163
zinc	Natural mineral: zinc	AR	159