

---

# **Quantitative Economics with Python using JAX**

**Thomas J. Sargent & John Stachurski**

**Aug 12, 2024**



# CONTENTS

<b>I</b>	<b>Introduction</b>	<b>3</b>
<b>1</b>	<b>About</b>	<b>5</b>
1.1	What is JAX? . . . . .	5
1.2	How to run these lectures . . . . .	6
1.3	Credits . . . . .	6
1.4	Prerequisites . . . . .	7
<b>2</b>	<b>An Introduction to JAX</b>	<b>9</b>
2.1	JAX as a NumPy Replacement . . . . .	10
2.2	Random Numbers . . . . .	14
2.3	JIT compilation . . . . .	15
2.4	Functional Programming . . . . .	18
2.5	Gradients . . . . .	19
2.6	Writing vectorized code . . . . .	20
2.7	Exercises . . . . .	23
<b>3</b>	<b>Adventures with Autodiff</b>	<b>25</b>
3.1	Overview . . . . .	25
3.2	What is automatic differentiation? . . . . .	26
3.3	Some experiments . . . . .	28
3.4	Gradient Descent . . . . .	33
3.5	Exercises . . . . .	38
<b>4</b>	<b>Newton's Method via JAX</b>	<b>41</b>
4.1	Overview . . . . .	41
4.2	Newton in one dimension . . . . .	42
4.3	An Equilibrium Problem . . . . .	44
4.4	Computation . . . . .	44
4.5	Exercises . . . . .	47
<b>II</b>	<b>Simulation</b>	<b>53</b>
<b>5</b>	<b>Inventory Dynamics</b>	<b>55</b>
5.1	Overview . . . . .	55
5.2	Sample paths . . . . .	56
5.3	Cross-sectional distributions . . . . .	57
5.4	Distribution dynamics . . . . .	61
5.5	Restock frequency . . . . .	63
<b>6</b>	<b>Kesten Processes and Firm Dynamics</b>	<b>67</b>

6.1	Overview . . . . .	67
6.2	Kesten processes . . . . .	68
<b>7</b>	<b>Wealth Distribution Dynamics</b>	<b>75</b>
7.1	Wealth dynamics . . . . .	76
7.2	Implementation . . . . .	77
7.3	Exercises . . . . .	87
<b>III</b>	<b>Asset Pricing</b>	<b>95</b>
<b>8</b>	<b>Asset Pricing: The Lucas Asset Pricing Model</b>	<b>97</b>
8.1	Overview . . . . .	97
8.2	The Lucas Model . . . . .	98
8.3	Computation . . . . .	102
8.4	Exercises . . . . .	108
<b>9</b>	<b>An Asset Pricing Problem</b>	<b>109</b>
9.1	Overview . . . . .	109
9.2	Pricing a single payoff . . . . .	110
9.3	Pricing a cash flow . . . . .	111
9.4	Choosing the stochastic discount factor . . . . .	112
9.5	Solving for the price-dividend ratio . . . . .	112
9.6	Code . . . . .	113
9.7	An Extended Example . . . . .	117
9.8	Numpy Version . . . . .	118
9.9	JAX Version . . . . .	123
9.10	A memory-efficient JAX version . . . . .	125
<b>IV</b>	<b>Dynamic Programming</b>	<b>129</b>
<b>10</b>	<b>Job Search</b>	<b>131</b>
10.1	Model . . . . .	131
10.2	Code . . . . .	132
10.3	Computing the solution . . . . .	134
10.4	Exercise . . . . .	136
<b>11</b>	<b>Optimal Savings I: Value Function Iteration</b>	<b>141</b>
11.1	Overview . . . . .	142
11.2	Starting with NumPy . . . . .	143
11.3	Switching to JAX . . . . .	146
11.4	Switching to vmap . . . . .	149
<b>12</b>	<b>Optimal Savings II: Alternative Algorithms</b>	<b>153</b>
12.1	Model primitives . . . . .	155
12.2	Operators . . . . .	156
12.3	Iteration . . . . .	158
12.4	Solvers . . . . .	159
12.5	Tests . . . . .	161
<b>13</b>	<b>Shortest Paths</b>	<b>169</b>
13.1	Overview . . . . .	169
13.2	Solving for Minimum Cost-to-Go . . . . .	170
13.3	Exercises . . . . .	172

<b>14 Optimal Investment</b>	<b>179</b>
<b>15 Inventory Management Model</b>	<b>193</b>
15.1 A model with constant discounting . . . . .	193
15.2 Time varying discount rates . . . . .	194
15.3 Numba implementation . . . . .	202
<b>16 Endogenous Grid Method</b>	<b>207</b>
16.1 Overview . . . . .	207
16.2 Setup . . . . .	208
16.3 Solution method . . . . .	209
16.4 Solutions . . . . .	214
<b>V Macroeconomic Models</b>	<b>219</b>
<b>17 Default Risk and Income Fluctuations</b>	<b>221</b>
17.1 Overview . . . . .	221
17.2 Structure . . . . .	223
17.3 Equilibrium . . . . .	225
17.4 Computation . . . . .	226
17.5 Results . . . . .	233
17.6 Exercises . . . . .	235
<b>18 The Aiyagari Model</b>	<b>241</b>
18.1 Overview . . . . .	241
18.2 Firms . . . . .	243
18.3 Households . . . . .	244
18.4 Equilibrium . . . . .	251
18.5 Exercises . . . . .	255
<b>19 The Hopenhayn Entry-Exit Model</b>	<b>259</b>
19.1 Outline . . . . .	259
19.2 The Model . . . . .	260
19.3 Code . . . . .	263
19.4 Solving the model . . . . .	268
19.5 Pareto tails . . . . .	272
19.6 Exercise . . . . .	273
<b>20 Bianchi Overborrowing Model</b>	<b>277</b>
20.1 Description of the model . . . . .	277
20.2 Markov dynamics . . . . .	279
20.3 Overborrowing Model . . . . .	282
20.4 Planner problem . . . . .	287
20.5 Numerical solution . . . . .	289
20.6 Exercise . . . . .	293
<b>VI Data and Empirics</b>	<b>295</b>
<b>21 Maximum Likelihood Estimation</b>	<b>297</b>
21.1 Overview . . . . .	297
21.2 MLE with numerical methods (JAX) . . . . .	298
21.3 MLE with <code>statsmodels</code> . . . . .	303

<b>VII Other</b>	<b>307</b>
<b>22 Troubleshooting</b>	<b>309</b>
22.1 Fixing Your Local Environment . . . . .	309
22.2 Reporting an Issue . . . . .	310
<b>23 References</b>	<b>311</b>
<b>24 Execution Statistics</b>	<b>313</b>
<b>Bibliography</b>	<b>315</b>
<b>Index</b>	<b>317</b>

This website presents a set of lectures on quantitative economic modeling using GPUs and [Google JAX](#).

- Introduction
  - *About*
  - *An Introduction to JAX*
  - *Adventures with Autodiff*
  - *Newton's Method via JAX*
- Simulation
  - *Inventory Dynamics*
  - *Kesten Processes and Firm Dynamics*
  - *Wealth Distribution Dynamics*
- Asset Pricing
  - *Asset Pricing: The Lucas Asset Pricing Model*
  - *An Asset Pricing Problem*
- Dynamic Programming
  - *Job Search*
  - *Optimal Savings I: Value Function Iteration*
  - *Optimal Savings II: Alternative Algorithms*
  - *Shortest Paths*
  - *Optimal Investment*
  - *Inventory Management Model*
  - *Endogenous Grid Method*
- Macroeconomic Models
  - *Default Risk and Income Fluctuations*
  - *The Aiyagari Model*
  - *The Hopenhayn Entry-Exit Model*
  - *Bianchi Overborrowing Model*
- Data and Empirics
  - *Maximum Likelihood Estimation*
- Other
  - *Troubleshooting*
  - *References*
  - *Execution Statistics*





## **Part I**

# **Introduction**



## ABOUT

Perhaps the single most notable feature of scientific computing in the past two decades is the rise and rise of parallel computation.

For example, the advanced artificial intelligence applications now shaking the worlds of business and academia require massive computer power to train, and the great majority of that computer power is supplied by GPUs.

For us economists, with our ever-growing need for more compute cycles, parallel computing provides both opportunities and new difficulties.

The main difficulty we face vis-a-vis parallel computation is accessibility.

Even for those with time to invest in careful parallelization of their programs, exploiting the full power of parallel hardware is challenging for non-experts.

Moreover, that hardware changes from year to year, so any human capital associated with mastering intricacies of a particular GPU has a very high depreciation rate.

For these reasons, we find [Google JAX](#) compelling.

In short, JAX makes high performance and parallel computing accessible (and fun!).

It provides a familiar array programming interface based on NumPy, and, as long as some simple conventions are adhered to, this code compiles to extremely efficient and well-parallelized machine code.

One of the most agreeable features of JAX is that the same code set and be run on either CPUs or GPUs, which allows users to test and develop locally, before deploying to a more powerful machine for heavier computations.

JAX is relatively easy to learn and highly portable, allowing us programmers to focus on the algorithms we want to implement, rather than particular features of our hardware.

This lecture series provides an introduction to using Google JAX for quantitative economics.

The rest of this page provides some background information on JAX, notes on how to run the lectures, and credits for our colleagues and RAs.

## 1.1 What is JAX?

JAX is an open source Python library developed by Google Research to support in-house artificial intelligence and machine learning.

JAX provides data types, functions and a compiler for fast linear algebra operations and automatic differentiation.

Loosely speaking, JAX is like [NumPy](#) with the addition of

- automatic differentiation
- automated GPU/TPU support

- a just-in-time compiler

In short, JAX delivers

1. high execution speeds on CPUs due to efficient parallelization and JIT compilation,
2. a powerful and convenient environment for GPU programming, and
3. the ability to efficiently differentiate smooth functions for optimization and estimation.

These features make JAX ideal for almost all quantitative economic modeling problems that require heavy-duty computing.

## 1.2 How to run these lectures

The easiest way to run these lectures is via [Google Colab](#).

JAX is pre-installed with GPU support on Colab and Colab provides GPU access even on the free tier.

Each lecture has a “play” button on the top right that you can use to launch the lecture on Colab.

You might also like to try using JAX locally.

If you do not own a GPU, you can still install JAX for the CPU by following the relevant [install instructions](#).

(We recommend that you install [Anaconda Python](#) first.)

If you do have a GPU, you can try installing JAX for the GPU by following the install instructions for GPUs.

(This is not always trivial but is starting to get easier.)

## 1.3 Credits

In building this lecture series, we had invaluable assistance from research assistants at QuantEcon and our QuantEcon colleagues.

In particular, we thank and credit

- [Shu Hu](#)
- [Smit Lunagariya](#)
- [Matthew McKay](#)
- [Humphrey Yang](#)
- [Hengcheng Zhang](#)
- [Frank Wu](#)

## 1.4 Prerequisites

We assume that readers have covered most of the QuantEcon lecture series [on Python programming](#).



## GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

This lecture provides a short introduction to [Google JAX](#).

Let's see if we have an active GPU:

```
!nvidia-smi
```

Mon Aug 12 03:58:50 2024

```

+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5 |
+-----+
+-----+
| GPU Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp  Perf              Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+=====+=====+=====+
|  0  Tesla T4                        On  | 00000001:00:00.0 Off |
| Off |
| N/A   42C    P0                     32W / 70W |      2MiB / 16384MiB |      0%
| Default |
|
| N/A |
+-----+
+-----+
+-----+
| Processes:
|

```

(continues on next page)

(continued from previous page)

```

| GPU  GI  CI      PID  Type  Process name                                GPU
↳Memory |                                         ↳
|      ID  ID                                         Usage↳
↳      |
|=====|
| No running processes found                      ↳
↳      |
+-----+
↳-----+

```

## 2.1 JAX as a NumPy Replacement

One way to use JAX is as a plug-in NumPy replacement. Let's look at the similarities and differences.

### 2.1.1 Similarities

The following import is standard, replacing `import numpy as np`:

```

import jax
import jax.numpy as jnp

```

Now we can use `jnp` in place of `np` for the usual array operations:

```
a = jnp.asarray((1.0, 3.2, -1.5))
```

```

2024-08-12 03:58:51.676938: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.

```

```
print(a)
```

```
[ 1.   3.2 -1.5]
```

```
print(jnp.sum(a))
```

```
2.6999998
```

```
print(jnp.mean(a))
```

```
0.9
```

```
print(jnp.dot(a, a))
```



```
13.490001
```

However, the array object `a` is not a NumPy array:

```
a
```

```
Array([ 1. ,  3.2, -1.5], dtype=float32)
```

```
type(a)
```

```
jaxlib.xla_extension.ArrayImpl
```

Even scalar-valued maps on arrays return JAX arrays.

```
jnp.sum(a)
```

```
Array(2.6999998, dtype=float32)
```

JAX arrays are also called “device arrays,” where term “device” refers to a hardware accelerator (GPU or TPU).

(In the terminology of GPUs, the “host” is the machine that launches GPU operations, while the “device” is the GPU itself.)

Operations on higher dimensional arrays are also similar to NumPy:

```
A = jnp.ones((2, 2))
B = jnp.identity(2)
A @ B
```

```
Array([[1., 1.],
       [1., 1.]], dtype=float32)
```

```
from jax.numpy import linalg
```

```
linalg.inv(B)  # Inverse of identity is identity
```

```
Array([[1., 0.],
       [0., 1.]], dtype=float32)
```

```
linalg.eigh(B)  # Computes eigenvalues and eigenvectors
```

```
EighResult(eigenvalues=Array([0.99999994, 0.99999994], dtype=float32),
           eigenvectors=Array([[1., 0.],
                               [0., 1.]], dtype=float32))
```

## 2.1.2 Differences

One difference between NumPy and JAX is that JAX currently uses 32 bit floats by default.

This is standard for GPU computing and can lead to significant speed gains with small loss of precision.

However, for some calculations precision matters. In these cases 64 bit floats can be enforced via the command

```
jax.config.update("jax_enable_x64", True)
```

Let's check this works:

```
jnp.ones(3)
```

```
Array([1., 1., 1.], dtype=float64)
```

As a NumPy replacement, a more significant difference is that arrays are treated as **immutable**.

For example, with NumPy we can write

```
import numpy as np
a = np.linspace(0, 1, 3)
a
```

```
array([0. , 0.5, 1. ])
```

and then mutate the data in memory:

```
a[0] = 1
a
```

```
array([1. , 0.5, 1. ])
```

In JAX this fails:

```
a = jnp.linspace(0, 1, 3)
a
```

```
Array([0. , 0.5, 1. ], dtype=float64)
```

```
a[0] = 1
```

```
-----
TypeError                                Traceback (most recent call last)
Cell In[20], line 1
----> 1 a[0] = 1

File /opt/conda/envs/quantecon/lib/python3.11/site-packages/jax/_src/numpy/array_
methods.py:286, in _unimplemented_setitem(self, i, x)
    281 def _unimplemented_setitem(self, i, x):
    282     msg = ("'{}' object does not support item assignment. JAX arrays are "
    283           "immutable. Instead of ``x[idx] = y``, use ``x = x.at[idx].
    ↪set(y)`` "
```

(continues on next page)

(continued from previous page)

```

284         "or another .at[] method: "
285         "https://jax.readthedocs.io/en/latest/_autosummary/jax.numpy.
↳ ndarray.at.html")
--> 286     raise TypeError(msg.format(type(self)))

TypeError: '<class 'jaxlib.xla_extension.ArrayImpl'>' object does not support item_
↳ assignment. JAX arrays are immutable. Instead of ``x[idx] = y``, use ``x = x.
↳ at[idx].set(y)`` or another .at[] method: https://jax.readthedocs.io/en/latest/_
↳ autosummary/jax.numpy.ndarray.at.html

```

In line with immutability, JAX does not support inplace operations:

```

a = np.array((2, 1))
a.sort()
a

```

```
array([1, 2])
```

```

a = jnp.array((2, 1))
a_new = a.sort()
a, a_new

```

```
(Array([2, 1], dtype=int64), Array([1, 2], dtype=int64))
```

The designers of JAX chose to make arrays immutable because JAX uses a functional programming style. More on this below.

Note that, while mutation is discouraged, it is in fact possible with `at`, as in

```

a = jnp.linspace(0, 1, 3)
id(a)

```

```
106127472
```

```
a
```

```
Array([0. , 0.5, 1. ], dtype=float64)
```

```
a.at[0].set(1)
```

```
Array([1. , 0.5, 1. ], dtype=float64)
```

We can check that the array is mutated by verifying its identity is unchanged:

```
id(a)
```

```
106127472
```

## 2.2 Random Numbers

Random numbers are also a bit different in JAX, relative to NumPy. Typically, in JAX, the state of the random number generator needs to be controlled explicitly.

```
import jax.random as random
```

First we produce a key, which seeds the random number generator.

```
key = random.PRNGKey(1)
```

```
type(key)
```

```
jaxlib.xla_extension.ArrayImpl
```

```
print(key)
```

```
[0 1]
```

Now we can use the key to generate some random numbers:

```
x = random.normal(key, (3, 3))  
x
```

```
Array([[ -1.35247421, -0.2712502 , -0.02920518],  
       [ 0.34706456,  0.5464053 , -1.52325812],  
       [ 0.41677264, -0.59710138, -0.5678208 ]], dtype=float64)
```

If we use the same key again, we initialize at the same seed, so the random numbers are the same:

```
random.normal(key, (3, 3))
```

```
Array([[ -1.35247421, -0.2712502 , -0.02920518],  
       [ 0.34706456,  0.5464053 , -1.52325812],  
       [ 0.41677264, -0.59710138, -0.5678208 ]], dtype=float64)
```

To produce a (quasi-) independent draw, best practice is to “split” the existing key:

```
key, subkey = random.split(key)
```

```
random.normal(key, (3, 3))
```

```
Array([[ 1.85374374, -0.37683949, -0.61276867],  
       [-1.91829718,  0.27219409,  0.54922246],  
       [ 0.40451442, -0.58726839, -0.63967753]], dtype=float64)
```

```
random.normal(subkey, (3, 3))
```

```
Array([[ -0.4300635 ,  0.22778552,  0.57241269],
       [-0.15969178,  0.46719192,  0.21165091],
       [ 0.84118631,  1.18671326, -0.16607783]], dtype=float64)
```

The function below produces  $k$  (quasi-) independent random  $n \times n$  matrices using this procedure.

```
def gen_random_matrices(key, n, k):
    matrices = []
    for _ in range(k):
        key, subkey = random.split(key)
        matrices.append(random.uniform(subkey, (n, n)))
    return matrices
```

```
matrices = gen_random_matrices(key, 2, 2)
for A in matrices:
    print(A)
```

```
[[0.97440813 0.3838544 ]
 [0.9790686  0.99981046]]
[[0.3473302  0.17157842]
 [0.89346686 0.01403153]]
```

One point to remember is that JAX expects tuples to describe array shapes, even for flat arrays. Hence, to get a one-dimensional array of normal random draws we use `(len, )` for the shape, as in

```
random.normal(key, (5, ))
```

```
Array([-0.64377279,  0.76961857, -0.29809604,  0.47858776, -2.00591299],
      dtype=float64)
```

## 2.3 JIT compilation

The JAX just-in-time (JIT) compiler accelerates logic within functions by fusing linear algebra operations into a single optimized kernel that the host can launch on the GPU / TPU (or CPU if no accelerator is detected).

### 2.3.1 A first example

To see the JIT compiler in action, consider the following function.

```
def f(x):
    a = 3*x + jnp.sin(x) + jnp.cos(x**2) - jnp.cos(2*x) - x**2 * 0.4 * x**1.5
    return jnp.sum(a)
```

Let's build an array to call the function on.

```
n = 50_000_000
x = jnp.ones(n)
```

How long does the function take to execute?

```
%time f(x).block_until_ready()
```

```
CPU times: user 341 ms, sys: 9.61 ms, total: 350 ms  
Wall time: 501 ms
```

```
Array(2.19896006e+08, dtype=float64)
```

---

**Note:** Here, in order to measure actual speed, we use the `block_until_ready()` method to hold the interpreter until the results of the computation are returned from the device. This is necessary because JAX uses asynchronous dispatch, which allows the Python interpreter to run ahead of GPU computations.

---

The code doesn't run as fast as we might hope, given that it's running on a GPU.

But if we run it a second time it becomes much faster:

```
%time f(x).block_until_ready()
```

```
CPU times: user 0 ns, sys: 2.15 ms, total: 2.15 ms  
Wall time: 147 ms
```

```
Array(2.19896006e+08, dtype=float64)
```

This is because the built in functions like `jnp.cos` are JIT compiled and the first run includes compile time.

Why would JAX want to JIT-compile built in functions like `jnp.cos` instead of just providing pre-compiled versions, like NumPy?

The reason is that the JIT compiler can specialize on the *size* of the array being used, which is helpful for parallelization.

For example, in running the code above, the JIT compiler produced a version of `jnp.cos` that is specialized to floating point arrays of size `n = 50_000_000`.

We can check this by calling `f` with a new array of different size.

```
m = 50_000_001  
y = jnp.ones(m)
```

```
%time f(y).block_until_ready()
```

```
CPU times: user 244 ms, sys: 8.25 ms, total: 253 ms  
Wall time: 401 ms
```

```
Array(2.19896011e+08, dtype=float64)
```

Notice that the execution time increases, because now new versions of the built-ins like `jnp.cos` are being compiled, specialized to the new array size.

If we run again, the code is dispatched to the correct compiled version and we get faster execution.

```
%time f(y).block_until_ready()
```

```
CPU times: user 2.53 ms, sys: 0 ns, total: 2.53 ms
Wall time: 102 ms
```

```
Array(2.19896011e+08, dtype=float64)
```

The compiled versions for the previous array size are still available in memory too, and the following call is dispatched to the correct compiled code.

```
%time f(x).block_until_ready()
```

```
CPU times: user 234 µs, sys: 2.02 ms, total: 2.25 ms
Wall time: 102 ms
```

```
Array(2.19896006e+08, dtype=float64)
```

## 2.3.2 Compiling the outer function

We can do even better if we manually JIT-compile the outer function.

```
f_jit = jax.jit(f)    # target for JIT compilation
```

Let's run once to compile it:

```
f_jit(x)
```

```
Array(2.19896006e+08, dtype=float64)
```

And now let's time it.

```
%time f_jit(x).block_until_ready()
```

```
CPU times: user 593 µs, sys: 0 ns, total: 593 µs
Wall time: 231 ms
```

```
Array(2.19896006e+08, dtype=float64)
```

Note the speed gain.

This is because the array operations are fused and no intermediate arrays are created.

Incidentally, a more common syntax when targetting a function for the JIT compiler is

```
@jax.jit
def f(x):
    a = 3*x + jnp.sin(x) + jnp.cos(x**2) - jnp.cos(2*x) - x**2 * 0.4 * x**1.5
    return jnp.sum(a)
```

## 2.4 Functional Programming

From JAX's documentation:

*When walking about the countryside of Italy, the people will not hesitate to tell you that JAX has “una anima di pura programmazione funzionale”.*

In other words, JAX assumes a functional programming style.

The major implication is that JAX functions should be pure.

A pure function will always return the same result if invoked with the same inputs.

In particular, a pure function has

- no dependence on global variables and
- no side effects

JAX will not usually throw errors when compiling impure functions but execution becomes unpredictable.

Here's an illustration of this fact, using global variables:

```
a = 1 # global

@jax.jit
def f(x):
    return a + x
```

```
x = jnp.ones(2)
```

```
f(x)
```

```
Array([2., 2.], dtype=float64)
```

In the code above, the global value `a=1` is fused into the jitted function.

Even if we change `a`, the output of `f` will not be affected — as long as the same compiled version is called.

```
a = 42
```

```
f(x)
```

```
Array([2., 2.], dtype=float64)
```

Changing the dimension of the input triggers a fresh compilation of the function, at which time the change in the value of `a` takes effect:

```
x = jnp.ones(3)
```

```
f(x)
```

```
Array([43., 43., 43.], dtype=float64)
```

Moral of the story: write pure functions when using JAX!



## 2.5 Gradients

JAX can use automatic differentiation to compute gradients.

This can be extremely useful for optimization and solving nonlinear systems.

We will see significant applications later in this lecture series.

For now, here's a very simple illustration involving the function

```
def f(x):
    return (x**2) / 2
```

Let's take the derivative:

```
f_prime = jax.grad(f)
```

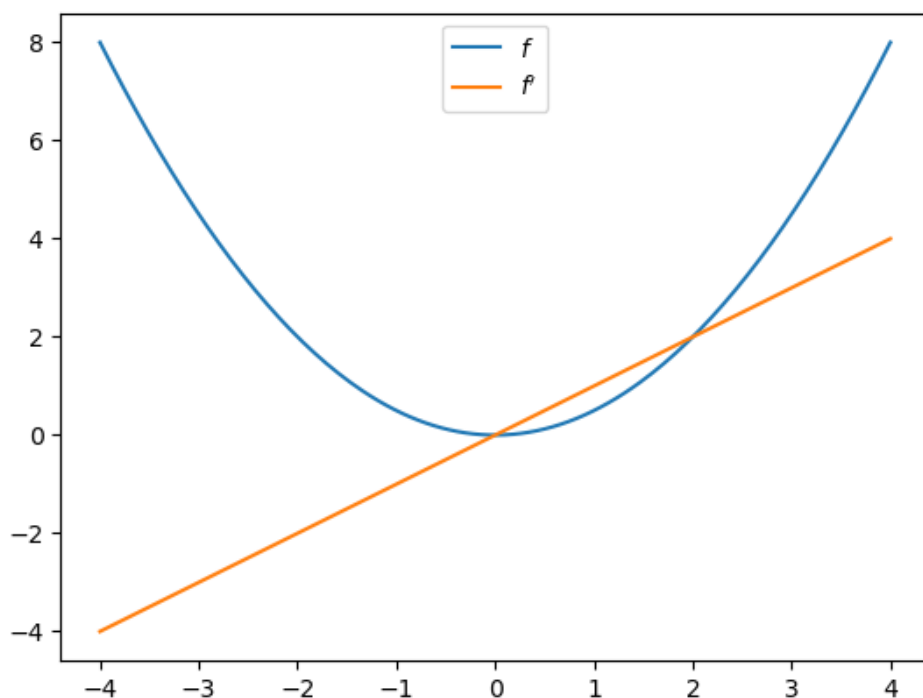
```
f_prime(10.0)
```

```
Array(10., dtype=float64, weak_type=True)
```

Let's plot the function and derivative, noting that  $f'(x) = x$ .

```
import matplotlib.pyplot as plt

fig, ax = plt.subplots()
x_grid = jnp.linspace(-4, 4, 200)
ax.plot(x_grid, f(x_grid), label="$f$")
ax.plot(x_grid, [f_prime(x) for x in x_grid], label="$f'$")
ax.legend(loc='upper center')
plt.show()
```



We defer further exploration of automatic differentiation with JAX until [Adventures with Autodiff](#).

## 2.6 Writing vectorized code

Writing fast JAX code requires shifting repetitive tasks from loops to array processing operations, so that the JAX compiler can easily understand the whole operation and generate more efficient machine code.

This procedure is called **vectorization** or **array programming**, and will be familiar to anyone who has used NumPy or MATLAB.

In most ways, vectorization is the same in JAX as it is in NumPy.

But there are also some differences, which we highlight here.

As a running example, consider the function

$$f(x, y) = \frac{\cos(x^2 + y^2)}{1 + x^2 + y^2}$$

Suppose that we want to evaluate this function on a square grid of  $x$  and  $y$  points and then plot it.

To clarify, here is the slow for loop version.

```
@jax.jit
def f(x, y):
    return jnp.cos(x**2 + y**2) / (1 + x**2 + y**2)

n = 80
x = jnp.linspace(-2, 2, n)
y = x

z_loops = np.empty((n, n))
```

```
%%time
for i in range(n):
    for j in range(n):
        z_loops[i, j] = f(x[i], y[j])
```

```
CPU times: user 4.33 s, sys: 1.23 s, total: 5.56 s
Wall time: 3.83 s
```

Even for this very small grid, the run time is extremely slow.

(Notice that we used a NumPy array for `z_loops` because we wanted to write to it.)

OK, so how can we do the same operation in vectorized form?

If you are new to vectorization, you might guess that we can simply write

```
z_bad = f(x, y)
```

But this gives us the wrong result because JAX doesn't understand the nested for loop.

```
z_bad.shape
```

```
(80,)
```

Here is what we actually wanted:

```
z_loops.shape
```

```
(80, 80)
```

To get the right shape and the correct nested for loop calculation, we can use a `meshgrid` operation designed for this purpose:

```
x_mesh, y_mesh = jnp.meshgrid(x, y)
```

Now we get what we want and the execution time is very fast.

```
%%time
z_mesh = f(x_mesh, y_mesh).block_until_ready()
```

```
CPU times: user 36.2 ms, sys: 591 µs, total: 36.8 ms
Wall time: 62.2 ms
```

Let's run again to eliminate compile time.

```
%%time
z_mesh = f(x_mesh, y_mesh).block_until_ready()
```

```
CPU times: user 426 µs, sys: 49 µs, total: 475 µs
Wall time: 270 µs
```

Let's confirm that we got the right answer.

```
jnp.allclose(z_mesh, z_loops)
```

```
Array(True, dtype=bool)
```

Now we can set up a serious grid and run the same calculation (on the larger grid) in a short amount of time.

```
n = 6000
x = jnp.linspace(-2, 2, n)
y = x
x_mesh, y_mesh = jnp.meshgrid(x, y)
```

```
%%time
z_mesh = f(x_mesh, y_mesh).block_until_ready()
```

```
CPU times: user 39.3 ms, sys: 0 ns, total: 39.3 ms
Wall time: 93.5 ms
```

Let's run again to get rid of compile time.

```
%%time
z_mesh = f(x_mesh, y_mesh).block_until_ready()
```

```
CPU times: user 496 µs, sys: 0 ns, total: 496 µs
Wall time: 27.4 ms
```

But there is one problem here: the mesh grids use a lot of memory.

```
x_mesh.nbytes + y_mesh.nbytes
```

```
576000000
```

By comparison, the flat array `x` is just

```
x.nbytes  # and y is just a pointer to x
```

```
48000
```

This extra memory usage can be a big problem in actual research calculations.

So let's try a different approach using `jax.vmap`

First we vectorize `f` in `y`.

```
f_vec_y = jax.vmap(f, in_axes=(None, 0))
```

In the line above, `(None, 0)` indicates that we are vectorizing in the second argument, which is `y`.

Next, we vectorize in the first argument, which is `x`.

```
f_vec = jax.vmap(f_vec_y, in_axes=(0, None))
```

With this construction, we can now call the function `f` on flat (low memory) arrays.

```
%%time
z_vmap = f_vec(x, y).block_until_ready()
```

```
CPU times: user 41.6 ms, sys: 0 ns, total: 41.6 ms
Wall time: 93.4 ms
```

We run it again to eliminate compile time.

```
%%time
z_vmap = f_vec(x, y).block_until_ready()
```

```
CPU times: user 1.14 ms, sys: 0 ns, total: 1.14 ms
Wall time: 25 ms
```

The execution time is essentially the same as the mesh operation but we are using much less memory.

And we produce the correct answer:

```
jnp.allclose(z_vmap, z_mesh)
```

```
Array(True, dtype=bool)
```

## 2.7 Exercises

### Exercise 2.7.1

In the Exercise section of [a lecture on Numba and parallelization](#), we used Monte Carlo to price a European call option. The code was accelerated by Numba-based multithreading.

Try writing a version of this operation for JAX, using all the same parameters.

If you are running your code on a GPU, you should be able to achieve significantly faster execution.

### Solution to Exercise 2.7.1

Here is one solution:

```
M = 10_000_000

n, β, K = 20, 0.99, 100
μ, ρ, v, S0, h0 = 0.0001, 0.1, 0.001, 10, 0

@jax.jit
def compute_call_price_jax(β=β,
                           μ=μ,
                           S0=S0,
                           h0=h0,
                           K=K,
                           n=n,
                           ρ=ρ,
                           v=v,
                           M=M,
                           key=jax.random.PRNGKey(1)):

    s = jnp.full(M, np.log(S0))
    h = jnp.full(M, h0)
    for t in range(n):
        key, subkey = jax.random.split(key)
        Z = jax.random.normal(subkey, (2, M))
        s = s + μ + jnp.exp(h) * Z[0, :]
        h = ρ * h + v * Z[1, :]
    expectation = jnp.mean(jnp.maximum(jnp.exp(s) - K, 0))

    return β**n * expectation
```

Let's run it once to compile it:

```
%%time
compute_call_price_jax().block_until_ready()
```

```
CPU times: user 11.4 s, sys: 275 ms, total: 11.7 s  
Wall time: 14 s
```

```
Array(180876.48840921, dtype=float64)
```

And now let's time it:

```
%%time  
compute_call_price_jax().block_until_ready()
```

```
CPU times: user 680 µs, sys: 0 ns, total: 680 µs  
Wall time: 459 ms
```

```
Array(180876.48840921, dtype=float64)
```

---

## ADVENTURES WITH AUTODIFF

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

### 3.1 Overview

This lecture gives a brief introduction to automatic differentiation using Google JAX.

Automatic differentiation is one of the key elements of modern machine learning and artificial intelligence.

As such it has attracted a great deal of investment and there are several powerful implementations available.

One of the best of these is the automatic differentiation routines contained in JAX.

While other software packages also offer this feature, the JAX version is particularly powerful because it integrates so well with other core components of JAX (e.g., JIT compilation and parallelization).

As we will see in later lectures, automatic differentiation can be used not only for AI but also for many problems faced in mathematical modeling, such as multi-dimensional nonlinear optimization and root-finding problems.

We need the following imports

```
import jax
import jax.numpy as jnp
import matplotlib.pyplot as plt
import numpy as np
```

Checking for a GPU:

```
!nvidia-smi
```

```
Mon Aug 12 03:21:50 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5                                    |
```

(continues on next page)

(continued from previous page)

```

|-----+-----+-----+
| GPU Name Persistence-M | Bus-Id Disp.A | Volatile Uncorr. |
| ECC |
| Fan Temp Perf Pwr:Usage/Cap | Memory-Usage | GPU-Util |
| Compute M. |
|
| MIG M. |
|=====|=====|=====|
| 0 Tesla T4 On | 00000001:00:00.0 Off | |
| Off |
| N/A 41C P0 32W / 70W | 2MiB / 16384MiB | 0% |
| Default |
|
| N/A |
|-----+-----+-----+
|
|-----+
|
|-----+
| Processes: |
| GPU GI CI PID Type Process name GPU |
| Memory |
| ID ID Usage |
|
|=====|=====|=====|
| No running processes found |
|
|-----+

```

### 3.2 What is automatic differentiation?

Autodiff is a technique for calculating derivatives on a computer.

### 3.2.1 Autodiff is not finite differences

The derivative of  $f(x) = \exp(2x)$  is

$$f'(x) = 2 \exp(2x)$$

A computer that doesn't know how to take derivatives might approximate this with the finite difference ratio

$$(Df)(x) := \frac{f(x+h) - f(x)}{h}$$

where  $h$  is a small positive number.

```
def f(x):  
    "Original function."  
    return np.exp(2 * x)
```

(continues on next page)

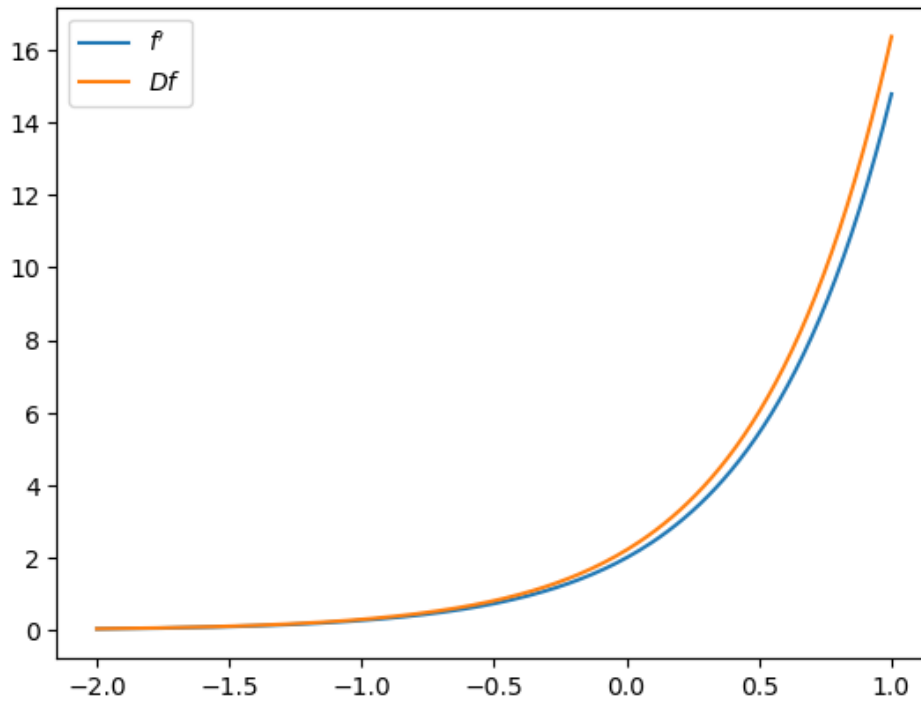


(continued from previous page)

```
def f_prime(x):
    "True derivative."
    return 2 * np.exp(2 * x)

def Df(x, h=0.1):
    "Approximate derivative (finite difference)."
    return (f(x + h) - f(x))/h

x_grid = np.linspace(-2, 1, 200)
fig, ax = plt.subplots()
ax.plot(x_grid, f_prime(x_grid), label="$f'$")
ax.plot(x_grid, Df(x_grid), label="$Df$")
ax.legend()
plt.show()
```



This kind of numerical derivative is often inaccurate and unstable.

One reason is that

$$\frac{f(x+h) - f(x)}{h} \approx \frac{0}{0}$$

Small numbers in the numerator and denominator causes rounding errors.

The situation is exponentially worse in high dimensions / with higher order derivatives

### 3.2.2 Autodiff is not symbolic calculus

Symbolic calculus tries to use rules for differentiation to produce a single closed-form expression representing a derivative.

```
from sympy import symbols, diff

m, a, b, x = symbols('m a b x')
f_x = (a*x + b)**m
f_x.diff((x, 6)) # 6-th order derivative
```

$$\frac{a^6 m (ax + b)^m (m^5 - 15m^4 + 85m^3 - 225m^2 + 274m - 120)}{(ax + b)^6}$$

Symbolic calculus is not well suited to high performance computing.

One disadvantage is that symbolic calculus cannot differentiate through control flow.

Also, using symbolic calculus might involve redundant calculations.

For example, consider

$$(fgh)' = (f'g + g'f)h + (fg)h'$$

If we evaluate at  $x$ , then we evaluate  $f(x)$  and  $g(x)$  twice each.

Also, computing  $f'(x)$  and  $f(x)$  might involve similar terms (e.g.,  $(f(x) = \exp(2x))' \implies f'(x) = 2f(x)$ ) but this is not exploited in symbolic algebra.

### 3.2.3 Autodiff

Autodiff produces functions that evaluates derivatives at numerical values passed in by the calling code, rather than producing a single symbolic expression representing the entire derivative.

Derivatives are constructed by breaking calculations into component parts via the chain rule.

The chain rule is applied until the point where the terms reduce to primitive functions that the program knows how to differentiate exactly (addition, subtraction, exponentiation, sine and cosine, etc.)

## 3.3 Some experiments

Let's start with some real-valued functions on  $\mathbb{R}$ .

### 3.3.1 A differentiable function

Let's test JAX's auto diff with a relatively simple function.

```
def f(x):
    return jnp.sin(x) - 2 * jnp.cos(3 * x) * jnp.exp(- x**2)
```

We use `grad` to compute the gradient of a real-valued function:

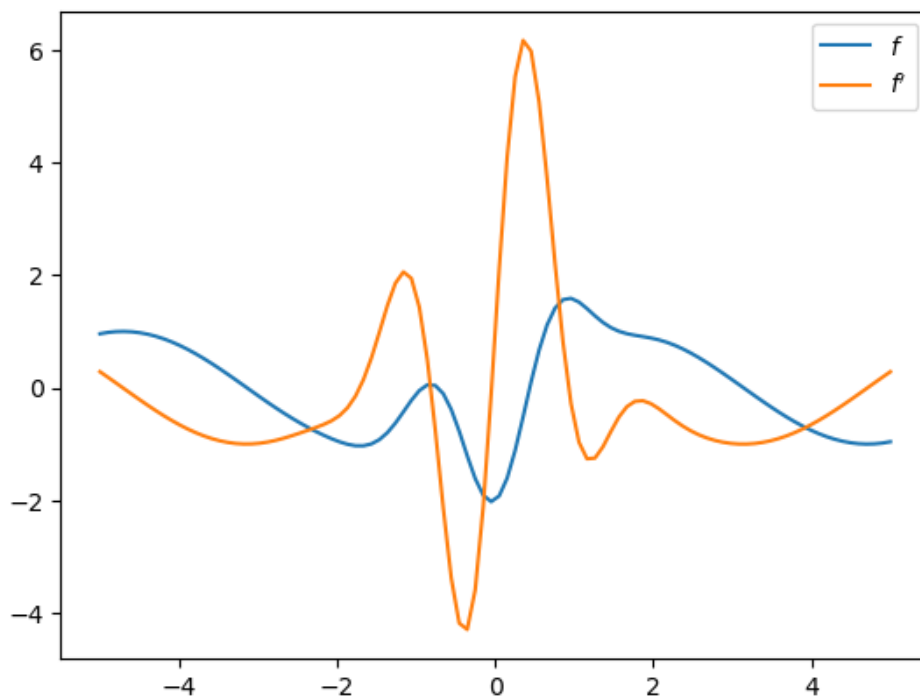
```
f_prime = jax.grad(f)
```

Let's plot the result:

```
x_grid = jnp.linspace(-5, 5, 100)
```

```
2024-08-12 03:21:51.108375: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```

```
fig, ax = plt.subplots()
ax.plot(x_grid, [f(x) for x in x_grid], label="$f$")
ax.plot(x_grid, [f_prime(x) for x in x_grid], label="$f'$")
ax.legend()
plt.show()
```



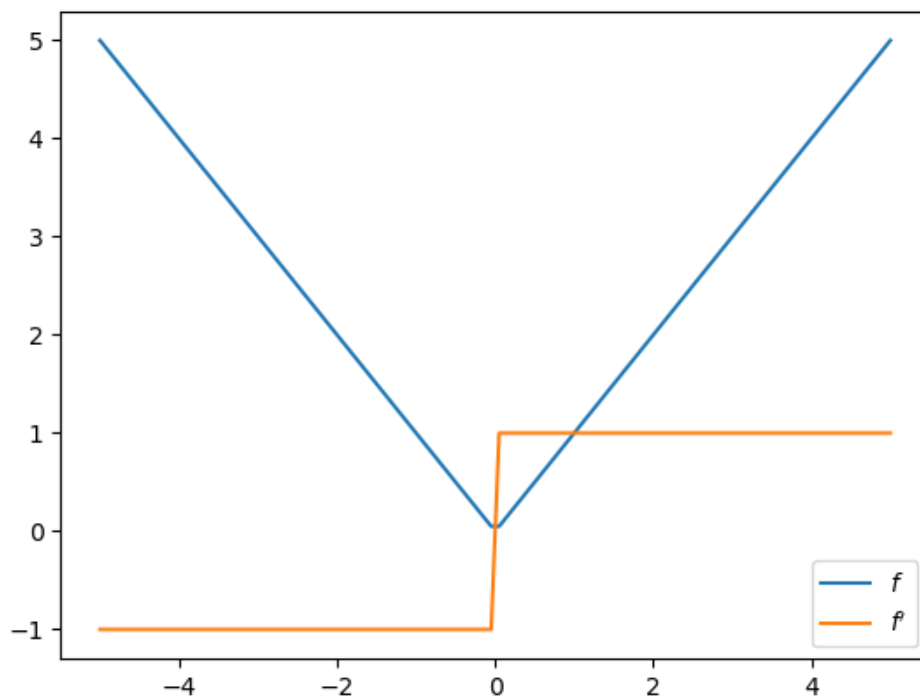
### 3.3.2 Absolute value function

What happens if the function is not differentiable?

```
def f(x):
    return jnp.abs(x)
```

```
f_prime = jax.grad(f)
```

```
fig, ax = plt.subplots()
ax.plot(x_grid, [f(x) for x in x_grid], label="$f$")
ax.plot(x_grid, [f_prime(x) for x in x_grid], label="$f'$")
ax.legend()
plt.show()
```



At the nondifferentiable point 0, `jax.grad` returns the right derivative:

```
f_prime(0.0)
```

```
Array(1., dtype=float32, weak_type=True)
```

### 3.3.3 Differentiating through control flow

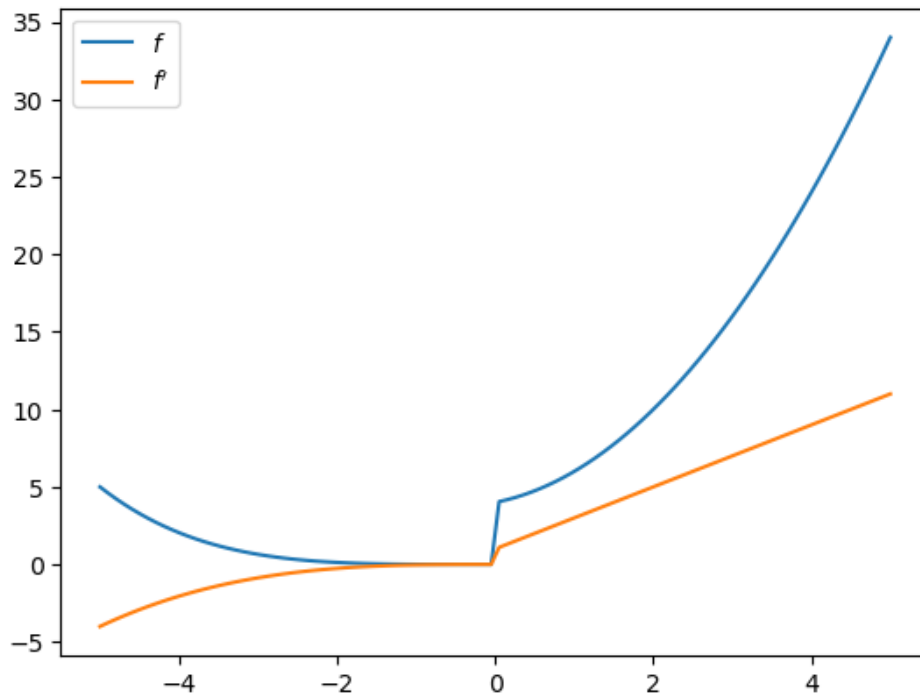
Let's try differentiating through some loops and conditions.

```
def f(x):
    def f1(x):
        for i in range(2):
            x *= 0.2 * x
        return x
    def f2(x):
        x = sum((x**i + i) for i in range(3))
        return x
    y = f1(x) if x < 0 else f2(x)
    return y
```

```
f_prime = jax.grad(f)
```

```
x_grid = jnp.linspace(-5, 5, 100)
```

```
fig, ax = plt.subplots()
ax.plot(x_grid, [f(x) for x in x_grid], label="$f$")
ax.plot(x_grid, [f_prime(x) for x in x_grid], label="$f'$")
ax.legend()
plt.show()
```

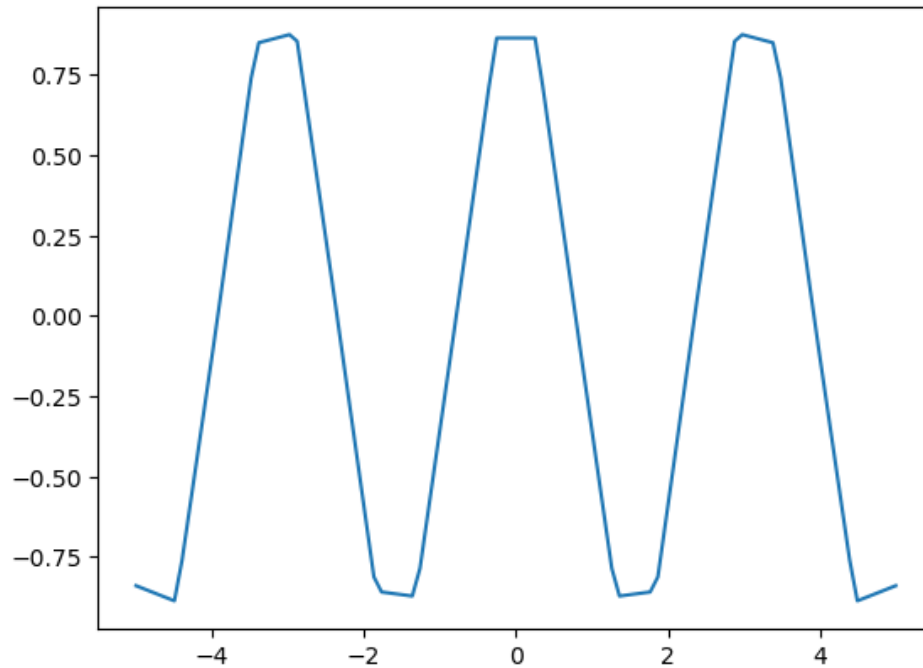


### 3.3.4 Differentiating through a linear interpolation

We can differentiate through linear interpolation, even though the function is not smooth:

```
n = 20
xp = jnp.linspace(-5, 5, n)
yp = jnp.cos(2 * xp)

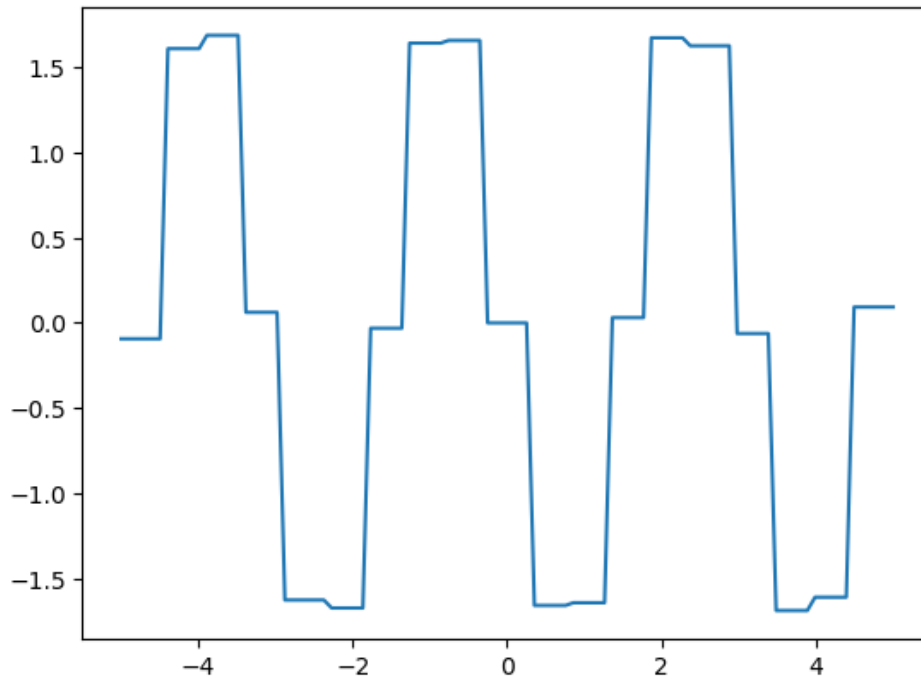
fig, ax = plt.subplots()
ax.plot(x_grid, jnp.interp(x_grid, xp, yp))
plt.show()
```



```
f_prime = jax.grad(jnp.interp)
```

```
f_prime_vec = jax.vmap(f_prime, in_axes=(0, None, None))
```

```
fig, ax = plt.subplots()
ax.plot(x_grid, f_prime_vec(x_grid, xp, yp))
plt.show()
```



## 3.4 Gradient Descent

Let's try implementing gradient descent.

As a simple application, we'll use gradient descent to solve for the OLS parameter estimates in simple linear regression.

### 3.4.1 A function for gradient descent

Here's an implementation of gradient descent.

```
def grad_descent(f,          # Function to be minimized
                 args,       # Extra arguments to the function
                 x0,          # Initial condition
                 λ=0.1,      # Initial learning rate
                 tol=1e-5,    # Tolerance
                 max_iter=1_000):
    """
    Minimize the function f via gradient descent, starting from guess x0.

    The learning rate is computed according to the Barzilai-Borwein method.

    """
    f_grad = jax.grad(f)
    x = jnp.array(x0)
    df = f_grad(x, args)
    ε = tol + 1
    i = 0
    while ε > tol and i < max_iter:
```

(continues on next page)

(continued from previous page)

```
new_x = x - λ * df
new_df = f_grad(new_x, args)
Δx = new_x - x
Δdf = new_df - df
λ = jnp.abs(Δx @ Δdf) / (Δdf @ Δdf)
ϵ = jnp.max(jnp.abs(Δx))
x, df = new_x, new_df
i += 1

return x
```

### 3.4.2 Simulated data

We're going to test our gradient descent function by minimizing a sum of least squares in a regression problem.

Let's generate some simulated data:

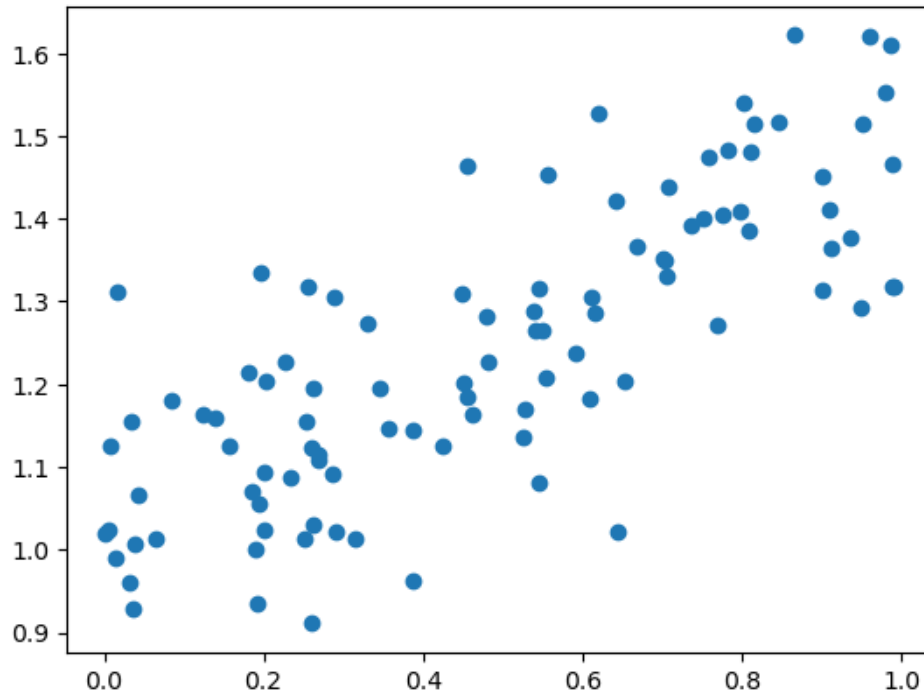
```
n = 100
key = jax.random.PRNGKey(1234)
x = jax.random.uniform(key, (n,))

α, β, σ = 0.5, 1.0, 0.1 # Set the true intercept and slope.
key, subkey = jax.random.split(key)
ϵ = jax.random.normal(subkey, (n,))

y = α * x + β + σ * ϵ
```

```
fig, ax = plt.subplots()
ax.scatter(x, y)
plt.show()
```





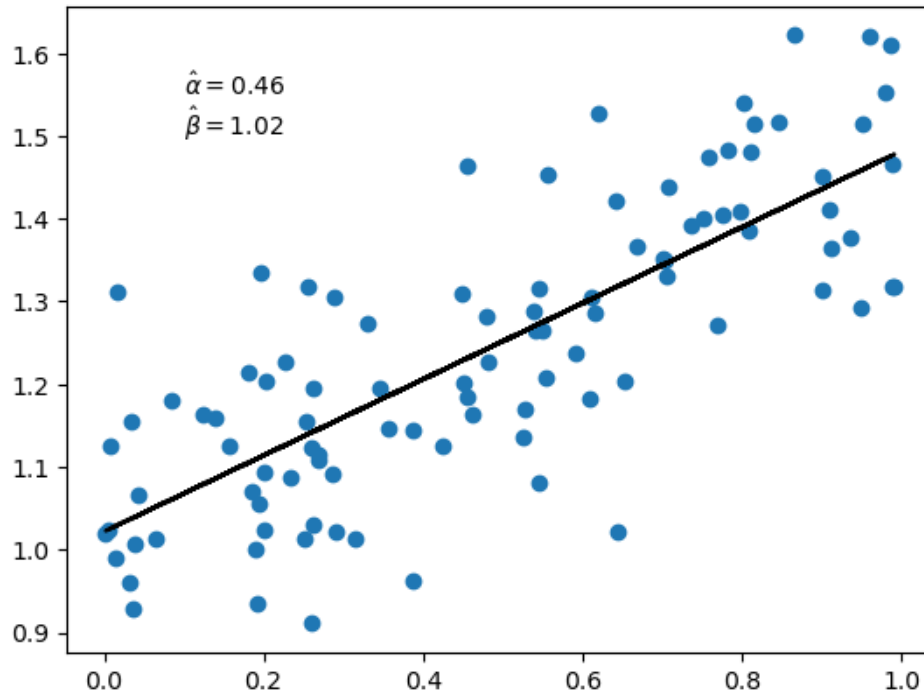
Let's start by calculating the estimated slope and intercept using closed form solutions.

```
mx = x.mean()
my = y.mean()
a_hat = jnp.sum((x - mx) * (y - my)) / jnp.sum((x - mx)**2)
beta_hat = my - a_hat * mx
```

```
a_hat, beta_hat
```

```
(Array(0.45994443, dtype=float32), Array(1.0219053, dtype=float32))
```

```
fig, ax = plt.subplots()
ax.scatter(x, y)
ax.plot(x, a_hat * x + beta_hat, 'k-')
ax.text(0.1, 1.55, rf'$\hat{\alpha} = {a_hat:.3}$')
ax.text(0.1, 1.50, rf'$\hat{\beta} = {beta_hat:.3}$')
plt.show()
```



### 3.4.3 Minimizing squared loss by gradient descent

Let's see if we can get the same values with our gradient descent function.

First we set up the least squares loss function.

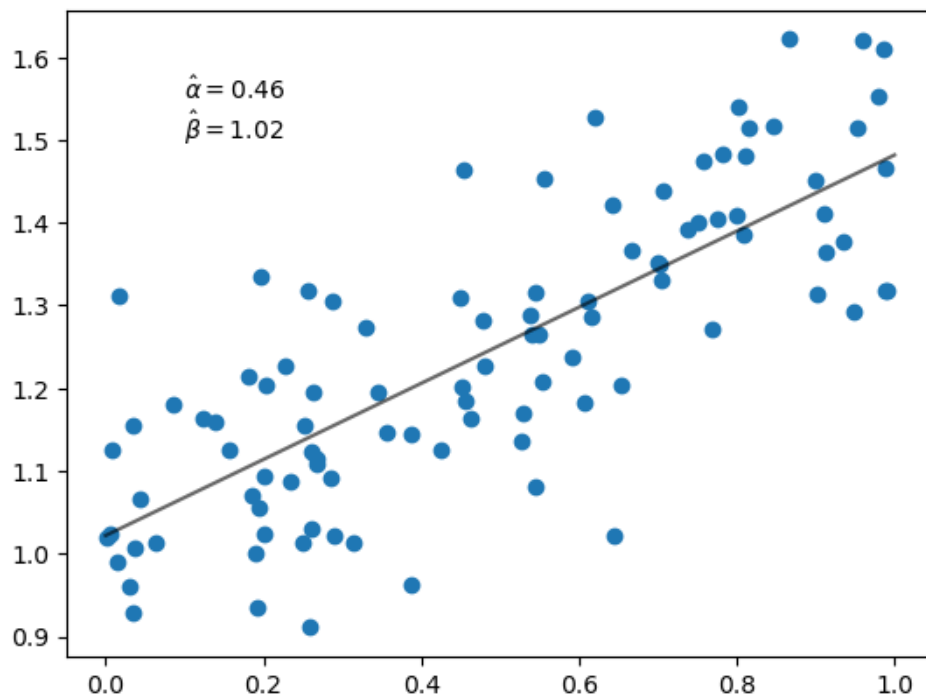
```
@jax.jit
def loss(params, data):
    a, b = params
    x, y = data
    return jnp.sum((y - a * x - b)**2)
```

Now we minimize it:

```
p0 = jnp.zeros(2)  # Initial guess for a, b
data = x, y
a_hat, b_hat = grad_descent(loss, data, p0)
```

Let's plot the results.

```
fig, ax = plt.subplots()
x_grid = jnp.linspace(0, 1, 100)
ax.scatter(x, y)
ax.plot(x_grid, a_hat * x_grid + b_hat, 'k-', alpha=0.6)
ax.text(0.1, 1.55, rf'$\hat{\alpha} = {a_hat:.3}$')
ax.text(0.1, 1.50, rf'$\hat{\beta} = {b_hat:.3}$')
plt.show()
```



Notice that we get the same estimates as we did from the closed form solutions.

### 3.4.4 Adding a squared term

Now let's try fitting a second order polynomial.

Here's our new loss function.

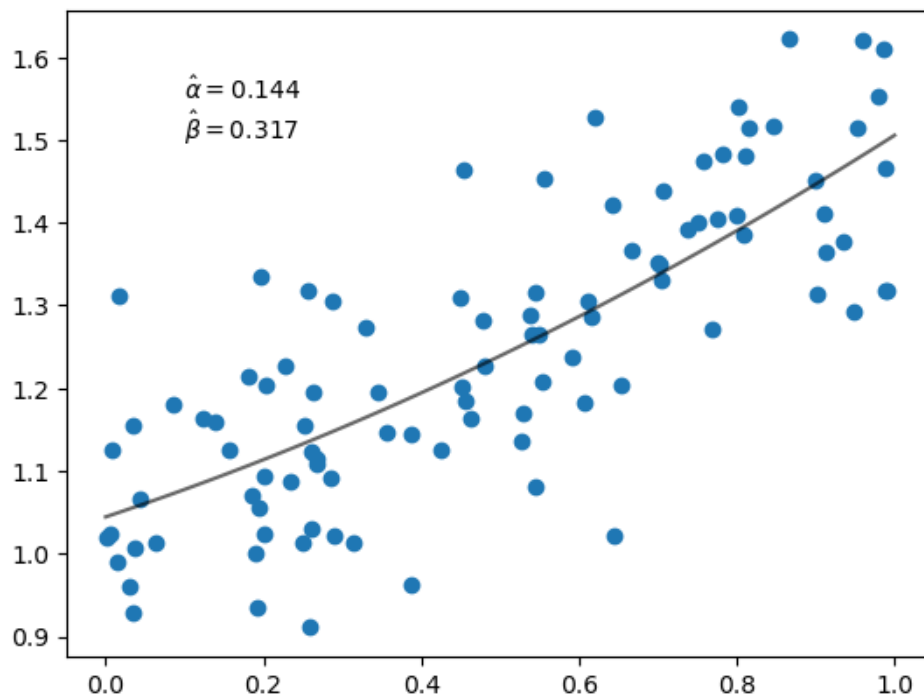
```
@jax.jit
def loss(params, data):
    a, b, c = params
    x, y = data
    return jnp.sum((y - a * x**2 - b * x - c)**2)
```

Now we're minimizing in three dimensions.

Let's try it.

```
p0 = jnp.zeros(3)
a_hat, b_hat, c_hat = grad_descent(loss, data, p0)

fig, ax = plt.subplots()
ax.scatter(x, y)
ax.plot(x_grid, a_hat * x_grid**2 + b_hat * x_grid + c_hat, 'k-', alpha=0.6)
ax.text(0.1, 1.55, rf'\hat \alpha = {a_hat:.3}')
ax.text(0.1, 1.50, rf'\hat \beta = {b_hat:.3}')
plt.show()
```



## 3.5 Exercises

---

### Exercise 3.5.1

The function `jnp.polyval` evaluates polynomials.

For example, if `len(p)` is 3, then `jnp.polyval(p, x)` returns

$$f(p, x) := p_0 x^2 + p_1 x + p_2$$

Use this function for polynomial regression.

The (empirical) loss becomes

$$\ell(p, x, y) = \sum_{i=1}^n (y_i - f(p, x_i))^2$$

Set  $k = 4$  and set the initial guess of `params` to `jnp.zeros(k)`.

Use gradient descent to find the array `params` that minimizes the loss function and plot the result (following the examples above).

---

### Solution to Exercise 3.5.1

Here's one solution.

```
def loss(params, data):  
    x, y = data  
    return jnp.sum((y - jnp.polyval(params, x))**2)
```

```

k = 4
p0 = jnp.zeros(k)
p_hat = grad_descent(loss, data, p0)
print('Estimated parameter vector:')
print(p_hat)
print('\n\n')

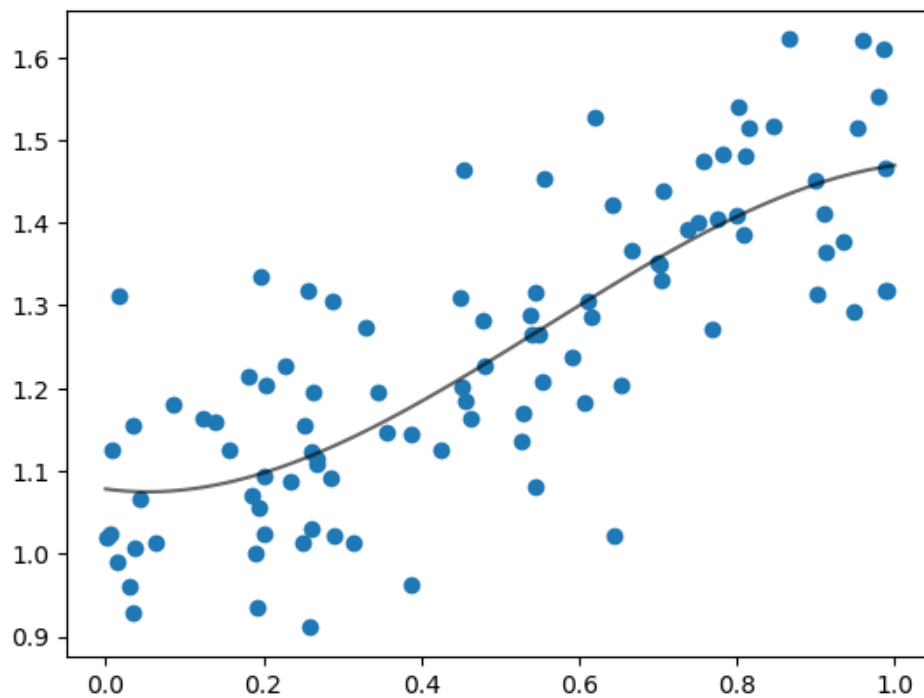
fig, ax = plt.subplots()
ax.scatter(x, y)
ax.plot(x_grid, jnp.polyval(p_hat, x_grid), 'k-', alpha=0.6)
plt.show()

```

```

Estimated parameter vector:
[-0.78131      1.3049686  -0.13246928  1.0781875 ]

```





## NEWTON'S METHOD VIA JAX

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

### 4.1 Overview

One of the key features of JAX is automatic differentiation.

We introduced this feature in *Adventures with Autodiff*.

In this lecture we apply automatic differentiation to the problem of computing economic equilibria via Newton's method.

Newton's method is a relatively simple root and fixed point solution algorithm, which we discussed in [a more elementary QuantEcon lecture](#).

JAX is ideally suited to implementing Newton's method efficiently, even in high dimensions.

We use the following imports in this lecture

```
import jax
import jax.numpy as jnp
from scipy.optimize import root
import matplotlib.pyplot as plt
```

Let's check the GPU we are running

```
!nvidia-smi
```

```
Mon Aug 12 04:00:24 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5 |
+-----+
| GPU   Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
(continues on next page)
```

(continued from previous page)

```

| Fan Temp Perf Pwr:Usage/Cap | Memory-Usage | GPU-Util |
└─Compute M. |
|
|
└─MIG M. |
|=====+=====+=====+
| 0 Tesla T4 On | 00000001:00:00.0 Off |
└─ Off |
| N/A 44C P0 32W / 70W | 2MiB / 16384MiB | 0% |
└─Default |
|
|
└─ N/A |
+-----+
└-----+
| Processes:
└─ |
| GPU GI CI PID Type Process name GPU
└─Memory |
| ID ID Usage
└─ |
|=====+=====+=====+
| No running processes found
└─ |
+-----+
└-----+

```

## 4.2 Newton in one dimension

As a warm up, let's implement Newton's method in JAX for a simple one-dimensional root-finding problem.

Let  $f$  be a function from  $\mathbb{R}$  to itself.

A **root** of  $f$  is an  $x \in \mathbb{R}$  such that  $f(x) = 0$ .

**Recall** that Newton's method for solving for the root of  $f$  involves iterating with the map  $q$  defined by

$$q(x) = x - \frac{f(x)}{f'(x)}$$

Here is a function called `newton` that takes a function  $f$  plus a scalar value  $x_0$ , iterates with  $q$  starting from  $x_0$ , and returns an approximate fixed point.

```
def newton(f, x_0, tol=1e-5):
    f_prime = jax.grad(f)
    def q(x):
        return x - f(x) / f_prime(x)

    error = tol + 1
    x = x_0
    while error > tol:
        y = q(x)
        error = abs(x - y)
```

(continues on next page)



(continued from previous page)

```

    x = y

    return x

```

The code above uses automatic differentiation to calculate  $f'$  via the call to `jax.grad`.

Let's test our newton routine on the function shown below.

```

f = lambda x: jnp.sin(4 * (x - 1/4)) + x + x**20 - 1
x = jnp.linspace(0, 1, 100)

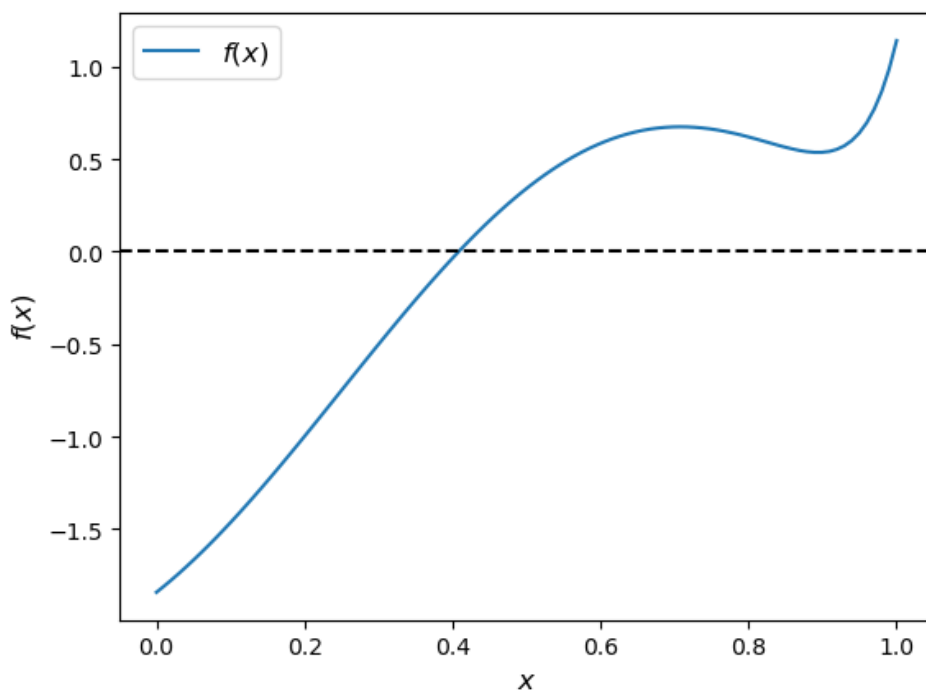
fig, ax = plt.subplots()
ax.plot(x, f(x), label='$f(x)$')
ax.axhline(ls='--', c='k')
ax.set_xlabel('$x$', fontsize=12)
ax.set_ylabel('$f(x)$', fontsize=12)
ax.legend(fontsize=12)
plt.show()

```

```

2024-08-12 04:00:25.172419: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.

```



Here we go

```

newton(f, 0.2)

```

```
Array(0.4082935, dtype=float32, weak_type=True)
```

This number looks to be close to the root, given the figure.

## 4.3 An Equilibrium Problem

Now let's move up to higher dimensions.

First we describe a market equilibrium problem we will solve with JAX via root-finding.

The market is for  $n$  goods.

(We are extending a two-good version of the market from [an earlier lecture](#).)

The supply function for the  $i$ -th good is

$$q_i^s(p) = b_i \sqrt{p_i}$$

which we write in vector form as

$$q^s(p) = b\sqrt{p}$$

(Here  $\sqrt{p}$  is the square root of each  $p_i$  and  $b\sqrt{p}$  is the vector formed by taking the pointwise product  $b_i \sqrt{p_i}$  at each  $i$ .)

The demand function is

$$q^d(p) = \exp(-Ap) + c$$

(Here  $A$  is an  $n \times n$  matrix containing parameters,  $c$  is an  $n \times 1$  vector and the  $\exp$  function acts pointwise (element-by-element) on the vector  $-Ap$ .)

The excess demand function is

$$e(p) = \exp(-Ap) + c - b\sqrt{p}$$

An **equilibrium price** vector is an  $n$ -vector  $p$  such that  $e(p) = 0$ .

The function below calculates the excess demand for given parameters

```
def e(p, A, b, c):  
    return jnp.exp(- A @ p) + c - b * jnp.sqrt(p)
```

## 4.4 Computation

In this section we describe and then implement the solution method.

### 4.4.1 Newton's Method

We use a multivariate version of Newton's method to compute the equilibrium price.

The rule for updating a guess  $p_n$  of the equilibrium price vector is

$$p_{n+1} = p_n - J_e(p_n)^{-1}e(p_n) \quad (4.1)$$

Here  $J_e(p_n)$  is the Jacobian of  $e$  evaluated at  $p_n$ .

Iteration starts from initial guess  $p_0$ .

Instead of coding the Jacobian by hand, we use automatic differentiation via `jax.jacobian()`.

```
def newton(f, x_0, tol=1e-5, max_iter=15):
    """
    A multivariate Newton root-finding routine.

    """
    x = x_0
    f_jac = jax.jacobian(f)
    @jax.jit
    def q(x):
        "Updates the current guess."
        return x - jnp.linalg.solve(f_jac(x), f(x))
    error = tol + 1
    n = 0
    while error > tol:
        n += 1
        if(n > max_iter):
            raise Exception('Max iteration reached without convergence')
        y = q(x)
        error = jnp.linalg.norm(x - y)
        x = y
        print(f'iteration {n}, error = {error}')
    return x
```

### 4.4.2 Application

Let's now apply the method just described to investigate a large market with 5,000 goods.

We randomly generate the matrix  $A$  and set the parameter vectors  $b, c$  to 1.

```
dim = 5_000
seed = 32

# Create a random matrix A and normalize the rows to sum to one
key = jax.random.PRNGKey(seed)
A = jax.random.uniform(key, [dim, dim])
s = jnp.sum(A, axis=0)
A = A / s

# Set up b and c
b = jnp.ones(dim)
c = jnp.ones(dim)
```

Here's our initial condition  $p_0$

```
init_p = jnp.ones(dim)
```

By combining the power of Newton's method, JAX accelerated linear algebra, automatic differentiation, and a GPU, we obtain a relatively small error for this high-dimensional problem in just a few seconds:

```
%%time
p = newton(lambda p: e(p, A, b, c), init_p).block_until_ready()
```

```
iteration 1, error = 29.977458953857422
iteration 2, error = 5.092828750610352
iteration 3, error = 0.10971637070178986
```

```
iteration 4, error = 5.167595372768119e-05
iteration 5, error = 1.1732883649528958e-05
iteration 6, error = 4.191016159893479e-06
CPU times: user 5.24 s, sys: 509 ms, total: 5.75 s
Wall time: 5.59 s
```

We run it again to eliminate compile time.

```
%%time
p = newton(lambda p: e(p, A, b, c), init_p).block_until_ready()
```

```
iteration 1, error = 29.977458953857422
iteration 2, error = 5.092828750610352
iteration 3, error = 0.10971637070178986
```

```
iteration 4, error = 5.167595372768119e-05
iteration 5, error = 1.1732883649528958e-05
iteration 6, error = 4.191016159893479e-06
CPU times: user 1.58 s, sys: 421 ms, total: 2 s
Wall time: 1.87 s
```

Here's the size of the error:

```
jnp.max(jnp.abs(e(p, A, b, c)))
```

```
Array(1.1920929e-07, dtype=float32)
```

With the same tolerance, SciPy's root function takes much longer to run, even with the Jacobian supplied.

```
%%time
solution = root(lambda p: e(p, A, b, c),
               init_p,
               jac=lambda p: jax.jacobian(e)(p, A, b, c),
               method='hybr',
               tol=1e-5)
```

```
CPU times: user 1min 51s, sys: 169 ms, total: 1min 51s
Wall time: 1min 51s
```

The result is also slightly less accurate:

```
p = solution.x
jnp.max(jnp.abs(e(p, A, b, c)))
```

```
Array(8.34465e-07, dtype=float32)
```

## 4.5 Exercises

### Exercise 4.5.1

Consider a three-dimensional extension of the [Solow fixed point problem](#) with

$$A = \begin{pmatrix} 2 & 3 & 3 \\ 2 & 4 & 2 \\ 1 & 5 & 1 \end{pmatrix}, \quad s = 0.2, \quad \alpha = 0.5, \quad \delta = 0.8$$

As before the law of motion is

$$k_{t+1} = g(k_t) \quad \text{where} \quad g(k) := sAk^\alpha + (1 - \delta)k$$

However  $k_t$  is now a  $3 \times 1$  vector.

Solve for the fixed point using Newton's method with the following initial values:

$$k1_0 = (1, 1, 1)$$

$$k2_0 = (3, 5, 5)$$

$$k3_0 = (50, 50, 50)$$

#### Hint:

- The computation of the fixed point is equivalent to computing  $k^*$  such that  $f(k^*) - k^* = 0$ .
- If you are unsure about your solution, you can start with the solved example:

$$A = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix}$$

with  $s = 0.3$ ,  $\alpha = 0.3$ , and  $\delta = 0.4$  and starting value:

$$k_0 = (1, 1, 1)$$

The result should converge to the [analytical solution](#).

### Solution to Exercise 4.5.1

Let's first define the parameters for this problem

```
A = jnp.array([[2.0, 3.0, 3.0],
               [2.0, 4.0, 2.0],
               [1.0, 5.0, 1.0]])

s = 0.2
α = 0.5
δ = 0.8
initIs = [jnp.ones(3),
          jnp.array([3.0, 5.0, 5.0]),
          jnp.repeat(50.0, 3)]
```

Then we define the multivariate version of the formula for the law of motion of capital

```
def multivariate_solow(k, A=A, s=s, α=α, δ=δ):
    return s * jnp.dot(A, k**α) + (1 - δ) * k
```

Let's run through each starting value and see the output

```
attempt = 1
for init in initIs:
    print(f'Attempt {attempt}: Starting value is {init} \n')
    %time k = newton(lambda k: multivariate_solow(k) - k, \
                    init).block_until_ready()
    print('-'*64)
    attempt += 1
```

Attempt 1: Starting value is [1. 1. 1.]

```
iteration 1, error = 50.496315002441406
iteration 2, error = 41.1093864440918
iteration 3, error = 4.294127464294434
iteration 4, error = 0.3854290544986725
iteration 5, error = 0.0054382034577429295
iteration 6, error = 8.92080606718082e-07
CPU times: user 138 ms, sys: 4.14 ms, total: 142 ms
Wall time: 177 ms
```

-----

Attempt 2: Starting value is [3. 5. 5.]

```
iteration 1, error = 2.0701100826263428
iteration 2, error = 0.12642373144626617
iteration 3, error = 0.0006017307168804109
iteration 4, error = 3.3717478231665154e-07
CPU times: user 77 ms, sys: 0 ns, total: 77 ms
Wall time: 71.2 ms
```

-----

Attempt 3: Starting value is [50. 50. 50.]

```
iteration 1, error = 73.00942993164062
```

```
iteration 2, error = 6.493789196014404
iteration 3, error = 0.6806989312171936
iteration 4, error = 0.016202213242650032
iteration 5, error = 1.0600916539260652e-05
```

(continues on next page)

(continued from previous page)

```
iteration 6, error = 9.830249609876773e-07
CPU times: user 173 ms, sys: 7.94 ms, total: 181 ms
Wall time: 171 ms
-----
```

We find that the results are invariant to the starting values.

But the number of iterations it takes to converge is dependent on the starting values.

Let substitute the output back into the formulate to check our last result

```
multivariate_solow(k) - k
```

```
Array([ 4.7683716e-07,  0.0000000e+00, -2.3841858e-07], dtype=float32)
```

Note the error is very small.

We can also test our results on the known solution

```
A = jnp.array([[2.0, 0.0, 0.0],
               [0.0, 2.0, 0.0],
               [0.0, 0.0, 2.0]])
s = 0.3
α = 0.3
δ = 0.4
init = jnp.repeat(1.0, 3)
%time k = newton(lambda k: multivariate_solow(k, A=A, s=s, α=α, δ=δ) - k, \
                 init).block_until_ready()
```

```
iteration 1, error = 1.5745922327041626
iteration 2, error = 0.21344946324825287
iteration 3, error = 0.002045975998044014
iteration 4, error = 8.259061701210157e-07
CPU times: user 144 ms, sys: 12.1 ms, total: 156 ms
Wall time: 170 ms
```

```
# Now we time it without compile
%time k = newton(lambda k: multivariate_solow(k, A=A, s=s, α=α, δ=δ) - k, \
                 init).block_until_ready()
```

```
iteration 1, error = 1.5745922327041626
```

```
iteration 2, error = 0.21344946324825287
iteration 3, error = 0.002045975998044014
iteration 4, error = 8.259061701210157e-07
CPU times: user 152 ms, sys: 4.06 ms, total: 156 ms
Wall time: 145 ms
```

The result is very close to the true solution but still slightly different.

We can increase the precision of the floating point numbers and restrict the tolerance to obtain a more accurate approximation (see detailed discussion in the [lecture on JAX](#))

```
# We will use 64 bit floats with JAX in order to increase the precision.
jax.config.update("jax_enable_x64", True)
init = init.astype('float64')

%time k = newton(lambda k: multivariate_solow(k, A=A, s=s, a=a, δ=δ) - k, \
                 init, tol=1e-7).block_until_ready()
```

```
iteration 1, error = 1.5745916432444333
iteration 2, error = 0.21344933091258958
iteration 3, error = 0.0020465547718452695
iteration 4, error = 2.0309190076799282e-07
iteration 5, error = 1.538370149106851e-15
CPU times: user 165 ms, sys: 8.08 ms, total: 173 ms
Wall time: 219 ms
```

```
# Now we time it without compile
%time k = newton(lambda k: multivariate_solow(k, A=A, s=s, a=a, δ=δ) - k, \
                 init, tol=1e-7).block_until_ready()
```

```
iteration 1, error = 1.5745916432444333
iteration 2, error = 0.21344933091258958
iteration 3, error = 0.0020465547718452695
iteration 4, error = 2.0309190076799282e-07
iteration 5, error = 1.538370149106851e-15
CPU times: user 79.9 ms, sys: 0 ns, total: 79.9 ms
Wall time: 97.9 ms
```

We can see it steps towards a more accurate solution.

---

### Exercise 4.5.2

In this exercise, let's try different initial values and check how Newton's method responds to different starting points.

Let's define a three-good problem with the following default values:

$$A = \begin{pmatrix} 0.2 & 0.1 & 0.7 \\ 0.3 & 0.2 & 0.5 \\ 0.1 & 0.8 & 0.1 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad \text{and} \quad c = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

For this exercise, use the following extreme price vectors as initial values:

$$\begin{aligned} p1_0 &= (5, 5, 5) \\ p2_0 &= (1, 1, 1) \\ p3_0 &= (4.5, 0.1, 4) \end{aligned}$$

Set the tolerance to  $10^{-15}$  for more accurate output.

---

**Hint:** Similar to exercise 1, enabling float64 for JAX can improve the precision of our results.

---

### Solution to Exercise 4.5.2

Define parameters and initial values



```
A = jnp.array([
    [0.2, 0.1, 0.7],
    [0.3, 0.2, 0.5],
    [0.1, 0.8, 0.1]
])
b = jnp.array([1.0, 1.0, 1.0])
c = jnp.array([1.0, 1.0, 1.0])
initIs = [jnp.repeat(5.0, 3),
          jnp.array([4.5, 0.1, 4.0])]
```

Let's run through each initial guess and check the output

```
attempt = 1
for init in initIs:
    print(f'Attempt {attempt}: Starting value is {init} \n')
    init = init.astype('float64')
    %time p = newton(lambda p: e(p, A, b, c), \
                    init, \
                    tol=1e-15, max_iter=15).block_until_ready()
    print('-'*64)
    attempt +=1
```

```
Attempt 1: Starting value is [5. 5. 5.]
```

```
iteration 1, error = 9.243805733085065
iteration 2, error = nan
CPU times: user 93.9 ms, sys: 24 µs, total: 93.9 ms
Wall time: 113 ms
-----
Attempt 2: Starting value is [4.5 0.1 4. ]

iteration 1, error = 4.892018895185869
iteration 2, error = 1.2120550201694784
iteration 3, error = 0.6942087122866175
iteration 4, error = 0.168951089180319
iteration 5, error = 0.005209730313222213
iteration 6, error = 4.3632751705775364e-06
iteration 7, error = 3.0460818773540415e-12
iteration 8, error = 0.0
CPU times: user 64.2 ms, sys: 4.06 ms, total: 68.3 ms
Wall time: 62.5 ms
-----
```

We can find that Newton's method may fail for some starting values.

Sometimes it may take a few initial guesses to achieve convergence.

Substitute the result back to the formula to check our result

```
e(p, A, b, c)
```

```
Array([0., 0., 0.], dtype=float64)
```

We can see the result is very accurate.



## **Part II**

# **Simulation**



## INVENTORY DYNAMICS

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

## 5.1 Overview

This lecture explores the inventory dynamics of a firm using so-called s-S inventory control.

Loosely speaking, this means that the firm

- waits until inventory falls below some value  $s$
- and then restocks with a bulk order of  $S$  units (or, in some models, restocks up to level  $S$ ).

We will be interested in the distribution of the associated Markov process, which can be thought of as cross-sectional distributions of inventory levels across a large number of firms, all of which

1. evolve independently and
2. have the same dynamics.

Note that we also studied this model in a [separate lecture](#), using Numba.

Here we study the same problem using JAX.

We will use the following imports:

```
import matplotlib.pyplot as plt
import numpy as np
import jax
import jax.numpy as jnp
from jax import random, lax
from collections import namedtuple
```

Here’s a description of our GPU:

```
!nvidia-smi
```

```

Mon Aug 12 03:26:11 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5 |
+-----+
+-----+
| GPU  Name                      Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp   Perf           Pwr:Usage/Cap |      Memory-Usage | GPU-Util  _
| Compute M. |
|
| MIG M. |
+=====+
|  0   Tesla T4                       On      | 00000001:00:00.0 Off  |
| Off |
| N/A   43C     P0              26W / 70W |      2MiB / 16384MiB |      0%    _
| Default |
|
| N/A |
+-----+
+-----+
+-----+
| Processes:
|
| GPU  GI    CI          PID    Type    Process name                        GPU
| Memory |
| ID     ID
|
+=====+
| No running processes found
|
+-----+

```

## 5.2 Sample paths

Consider a firm with inventory  $X_t$ .

The firm waits until  $X_t \leq s$  and then restocks up to  $S$  units.

It faces stochastic demand  $\{D_t\}$ , which we assume is IID across time and firms.

With notation  $a^+ := \max\{a, 0\}$ , inventory dynamics can be written as

$$X_{t+1} = \begin{cases} (S - D_{t+1})^+ & \text{if } X_t \leq s \\ (X_t - D_{t+1})^+ & \text{if } X_t > s \end{cases}$$

In what follows, we will assume that each  $D_t$  is lognormal, so that

$$D_t = \exp(\mu + \sigma Z_t)$$

where  $\mu$  and  $\sigma$  are parameters and  $\{Z_t\}$  is IID and standard normal.

Here's a `namedtuple` that stores parameters.

```
Parameters = namedtuple('Parameters', ['s', 'S', 'μ', 'σ'])

# Create a default instance
params = Parameters(s=10, S=100, μ=1.0, σ=0.5)
```

## 5.3 Cross-sectional distributions

Now let's look at the marginal distribution  $\psi_T$  of  $X_T$  for some fixed  $T$ .

The probability distribution  $\psi_T$  is the time  $T$  distribution of firm inventory levels implied by the model.

We will approximate this distribution by

1. fixing  $n$  to be some large number, indicating the number of firms in the simulation,
2. fixing  $T$ , the time period we are interested in,
3. generating  $n$  independent draws from some fixed distribution  $\psi_0$  that gives the initial cross-section of inventories for the  $n$  firms, and
4. shifting this distribution forward in time  $T$  periods, updating each firm  $T$  times via the dynamics described above (independent of other firms).

We will then visualize  $\psi_T$  by histogramming the cross-section.

We will use the following code to update the cross-section of firms by one period.

```
@jax.jit
def update_cross_section(params, X_vec, D):
    """
    Update by one period a cross-section of firms with inventory levels given by
    X_vec, given the vector of demand shocks in D.

    * D[i] is the demand shock for firm i with current inventory X_vec[i]

    """
    # Unpack
    s, S = params.s, params.S
    # Restock if the inventory is below the threshold
    X_new = jnp.where(X_vec <= s,
                     jnp.maximum(S - D, 0), jnp.maximum(X_vec - D, 0))
    return X_new
```

### 5.3.1 For loop version

Now we provide code to compute the cross-sectional distribution  $\psi_T$  given some initial distribution  $\psi_0$  and a positive integer  $T$ .

In this code we use an ordinary Python `for` loop to step forward through time

While Python loops are slow, this approach is reasonable here because efficiency of outer loops has far less influence on runtime than efficiency of inner loops.

(Below we will squeeze out more speed by compiling the outer loop as well as the update rule.)

In the code below, the initial distribution  $\psi_0$  takes all firms to have initial inventory `x_init`.

```
def compute_cross_section(params, x_init, T, key, num_firms=50_000):
    # Set up initial distribution
    X_vec = jnp.full((num_firms, ), x_init)
    # Loop
    for i in range(T):
        Z = random.normal(key, shape=(num_firms, ))
        D = jnp.exp(params.μ + params.σ * Z)

        X_vec = update_cross_section(params, X_vec, D)
        _, key = random.split(key)

    return X_vec
```

We'll use the following specification

```
x_init = 50
T = 500
# Initialize random number generator
key = random.PRNGKey(10)
```

```
2024-08-12 03:26:11.944704: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```

Let's look at the timing.

```
%time X_vec = compute_cross_section(params, \
    x_init, T, key).block_until_ready()
```

```
CPU times: user 813 ms, sys: 156 ms, total: 969 ms
Wall time: 801 ms
```

Let's run again to eliminate compile time.

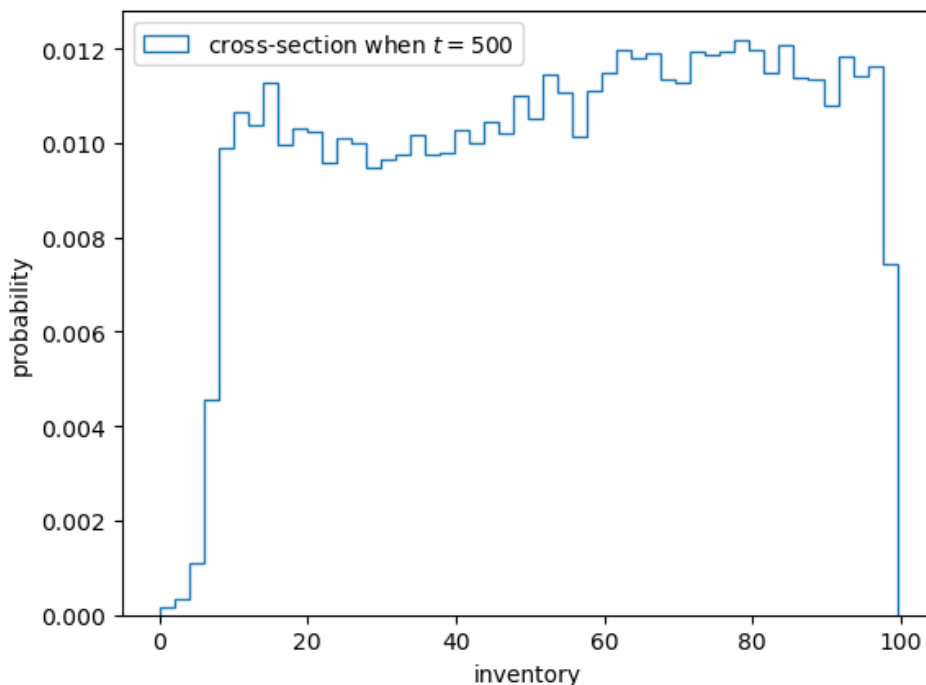
```
%time X_vec = compute_cross_section(params, \
    x_init, T, key).block_until_ready()
```

```
CPU times: user 318 ms, sys: 107 ms, total: 425 ms
Wall time: 244 ms
```

Here's a histogram of inventory levels at time  $T$ .

```
fig, ax = plt.subplots()
ax.hist(X_vec, bins=50,
        density=True,
        histtype='step',
        label=f'cross-section when $t = {T}$')
ax.set_xlabel('inventory')
ax.set_ylabel('probability')
ax.legend()
plt.show()
```





### 5.3.2 Compiling the outer loop

Now let's see if we can gain some speed by compiling the outer loop, which steps through the time dimension.

We will do this using `jax.jit` and a `fori_loop`, which is a compiler-ready version of a `for` loop provided by JAX.

```
def compute_cross_section_fori(params, x_init, T, key, num_firms=50_000):

    s, S, μ, σ = params.s, params.S, params.μ, params.σ
    X = jnp.full((num_firms, ), x_init)

    # Define the function for each update
    def fori_update(t, inputs):
        # Unpack
        X, key = inputs
        # Draw shocks using key
        Z = random.normal(key, shape=(num_firms,))
        D = jnp.exp(μ + σ * Z)
        # Update X
        X = jnp.where(X <= s,
                      jnp.maximum(S - D, 0),
                      jnp.maximum(X - D, 0))
        # Refresh the key
        key, subkey = random.split(key)
        return X, subkey

    # Loop t from 0 to T, applying fori_update each time.
    # The initial condition for fori_update is (X, key).
    X, key = lax.fori_loop(0, T, fori_update, (X, key))

    return X
```

(continues on next page)

(continued from previous page)

```
# Compile taking T and num_firms as static (changes trigger recompile)
compute_cross_section_for1 = jax.jit(
    compute_cross_section_for1, static_argnums=(2, 4))
```

Let's see how fast this runs with compile time.

```
%time X_vec = compute_cross_section_for1(params, \
    x_init, T, key).block_until_ready()
```

```
CPU times: user 459 ms, sys: 25.2 ms, total: 484 ms
Wall time: 424 ms
```

And let's see how fast it runs without compile time.

```
%time X_vec = compute_cross_section_for1(params, \
    x_init, T, key).block_until_ready()
```

```
CPU times: user 1.82 ms, sys: 461 µs, total: 2.28 ms
Wall time: 7.99 ms
```

Compared to the original version with a pure Python outer loop, we have produced a nontrivial speed gain.

This is due to the fact that we have compiled the whole operation.

### 5.3.3 Further vectorization

For relatively small problems, we can make this code run even faster by generating all random variables at once.

This improves efficiency because we are taking more operations out of the loop.

```
def compute_cross_section_for1(params, x_init, T, key, num_firms=50_000):

    s, S, μ, σ = params.s, params.S, params.μ, params.σ
    X = jnp.full((num_firms, ), x_init)
    Z = random.normal(key, shape=(T, num_firms))
    D = jnp.exp(μ + σ * Z)

    def update_cross_section(i, X):
        X = jnp.where(X <= s,
                      jnp.maximum(S - D[i, :], 0),
                      jnp.maximum(X - D[i, :], 0))
        return X

    X = lax.fori_loop(0, T, update_cross_section, X)

    return X

# Compile taking T and num_firms as static (changes trigger recompile)
compute_cross_section_for1 = jax.jit(
    compute_cross_section_for1, static_argnums=(2, 4))
```

Let's test it with compile time included.

```
%time X_vec = compute_cross_section_fori(params, \
                                         x_init, T, key).block_until_ready()
```

```
CPU times: user 260 ms, sys: 0 ns, total: 260 ms
Wall time: 243 ms
```

Let's run again to eliminate compile time.

```
%time X_vec = compute_cross_section_fori(params, \
                                         x_init, T, key).block_until_ready()
```

```
CPU times: user 3.76 ms, sys: 0 ns, total: 3.76 ms
Wall time: 6.15 ms
```

On one hand, this version is faster than the previous one, where random variables were generated inside the loop.

On the other hand, this implementation consumes far more memory, as we need to store large arrays of random draws.

The high memory consumption becomes problematic for large problems.

## 5.4 Distribution dynamics

Next let's take a look at how the distribution sequence evolves over time.

We will go back to using ordinary Python `for` loops.

Here is code that repeatedly shifts the cross-section forward while recording the cross-section at the dates in `sample_dates`.

```
def shift_forward_and_sample(x_init, params, sample_dates,
                             key, num_firms=50_000, sim_length=750):

    X = res = jnp.full((num_firms, ), x_init)

    # Use for loop to update X and collect samples
    for i in range(sim_length):
        Z = random.normal(key, shape=(num_firms, ))
        D = jnp.exp(params.μ + params.σ * Z)

        X = update_cross_section(params, X, D)
        _, key = random.split(key)

    # draw a sample at the sample dates
    if (i+1 in sample_dates):
        res = jnp.vstack((res, X))

    return res[1:]
```

Let's test it

```
x_init = 50
num_firms = 10_000
sample_dates = 10, 50, 250, 500, 750
```

(continues on next page)

(continued from previous page)

```
key = random.PRNGKey(10)

%time X = shift_forward_and_sample(x_init, params, \
                                   sample_dates, key).block_until_ready()
```

```
CPU times: user 642 ms, sys: 204 ms, total: 846 ms
Wall time: 639 ms
```

We run the code again to eliminate compile time.

```
%time X = shift_forward_and_sample(x_init, params, \
                                   sample_dates, key).block_until_ready()
```

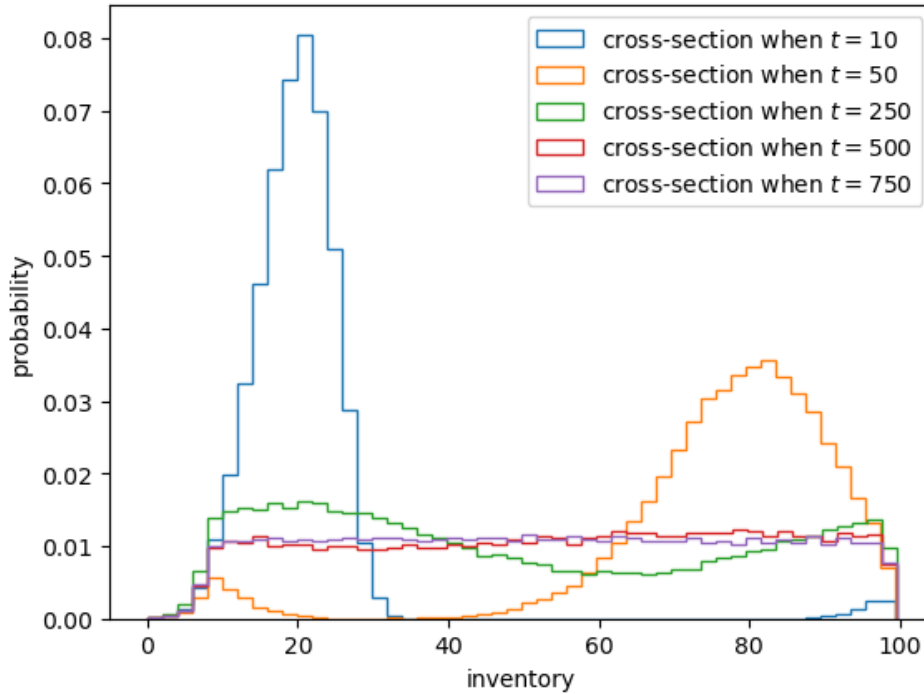
```
CPU times: user 476 ms, sys: 179 ms, total: 656 ms
Wall time: 390 ms
```

Let's plot the output.

```
fig, ax = plt.subplots()

for i, date in enumerate(sample_dates):
    ax.hist(X[i, :], bins=50,
            density=True,
            histtype='step',
            label=f'cross-section when $t = {date}$')

ax.set_xlabel('inventory')
ax.set_ylabel('probability')
ax.legend()
plt.show()
```



This model for inventory dynamics is asymptotically stationary, with a unique stationary distribution.

In particular, the sequence of marginal distributions  $\{\psi_t\}$  converges to a unique limiting distribution that does not depend on initial conditions.

Although we will not prove this here, we can see it in the simulation above.

By  $t = 500$  or  $t = 750$  the distributions are barely changing.

If you test a few different initial conditions, you will see that they do not affect long-run outcomes.

## 5.5 Restock frequency

As an exercise, let's study the probability that firms need to restock over a given time period.

In the exercise, we will

- set the starting stock level to  $X_0 = 70$  and
- calculate the proportion of firms that need to order twice or more in the first 50 periods.

This proportion approximates the probability of the event when the sample size is large.

### 5.5.1 For loop version

We start with an easier `for` loop implementation

```
# Define a jitted function for each update
@jax.jit
def update_stock(n_restock, X, params, D):
    n_restock = jnp.where(X <= params.s,
                          n_restock + 1,
                          n_restock)
    X = jnp.where(X <= params.s,
                  jnp.maximum(params.S - D, 0),
                  jnp.maximum(X - D, 0))
    return n_restock, X, key

def compute_freq(params, key,
                 x_init=70,
                 sim_length=50,
                 num_firms=1_000_000):

    # Prepare initial arrays
    X = jnp.full((num_firms, ), x_init)

    # Stack the restock counter on top of the inventory
    n_restock = jnp.zeros((num_firms, ))

    # Use a for loop to perform the calculations on all states
    for i in range(sim_length):
        Z = random.normal(key, shape=(num_firms, ))
        D = jnp.exp(params.μ + params.σ * Z)
        n_restock, X, key = update_stock(
            n_restock, X, params, D)
        key = random.fold_in(key, i)

    return jnp.mean(n_restock > 1, axis=0)
```

```
key = random.PRNGKey(27)

%time freq = compute_freq(params, key).block_until_ready()
```

```
CPU times: user 661 ms, sys: 0 ns, total: 661 ms
Wall time: 722 ms
```

We run the code again to get rid of compile time.

```
%time freq = compute_freq(params, key).block_until_ready()
```

```
CPU times: user 60 ms, sys: 7.73 ms, total: 67.7 ms
Wall time: 48.3 ms
```

```
print(f"Frequency of at least two stock outs = {freq}")
```

```
Frequency of at least two stock outs = 0.4472379982471466
```

**Exercise 5.5.1**

Write a `fori_loop` version of the last function. See if you can increase the speed while generating a similar answer.

**Solution to Exercise 5.5.1**

Here is a `lax.fori_loop` version that JIT compiles the whole function

```
@jax.jit
def compute_freq(params, key,
                 x_init=70,
                 sim_length=50,
                 num_firms=1_000_000):

    s, S,  $\mu$ ,  $\sigma$  = params.s, params.S, params. $\mu$ , params. $\sigma$ 
    # Prepare initial arrays
    X = jnp.full((num_firms, ), x_init)
    Z = random.normal(key, shape=(sim_length, num_firms))
    D = jnp.exp( $\mu$  +  $\sigma$  * Z)

    # Stack the restock counter on top of the inventory
    restock_count = jnp.zeros((num_firms, ))
    Xs = (X, restock_count)

    # Define the function for each update
    def update_cross_section(i, Xs):
        # Separate the inventory and restock counter
        x, restock_count = Xs[0], Xs[1]
        restock_count = jnp.where(x <= s,
                                restock_count + 1,
                                restock_count)

        x = jnp.where(x <= s,
                     jnp.maximum(S - D[i], 0),
                     jnp.maximum(x - D[i], 0))

        Xs = (x, restock_count)
        return Xs

    # Use lax.fori_loop to perform the calculations on all states
    X_final = lax.fori_loop(0, sim_length, update_cross_section, Xs)

    return jnp.mean(X_final[1] > 1)
```

Note the time the routine takes to run, as well as the output

```
%time freq = compute_freq(params, key).block_until_ready()
```

```
CPU times: user 296 ms, sys: 0 ns, total: 296 ms
Wall time: 294 ms
```

We run the code again to eliminate the compile time.

```
%time freq = compute_freq(params, key).block_until_ready()
```

```
CPU times: user 1.28 ms, sys: 0 ns, total: 1.28 ms  
Wall time: 7.12 ms
```

```
print(f"Frequency of at least two stock outs = {freq}")
```

```
Frequency of at least two stock outs = 0.44674399495124817
```

---



## KESTEN PROCESSES AND FIRM DYNAMICS

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

In addition to JAX and Anaconda, this lecture will need the following libraries:

```
!pip install quantecon
```

### 6.1 Overview

This lecture describes Kesten processes, which are an important class of stochastic processes, and an application of firm dynamics.

The lecture draws on [an earlier QuantEcon lecture](#), which uses Numba to accelerate the computations.

In that earlier lecture you can find a more detailed discussion of the concepts involved.

This lecture focuses on implementing the same computations in JAX.

Let’s start with some imports:

```
import matplotlib.pyplot as plt
import quantecon as qe
import jax
import jax.numpy as jnp
from jax import random
from jax import lax
```

Let’s check the GPU we are running

```
!nvidia-smi
```

```
Mon Aug 12 03:59:32 2024
```

```
+-----+
| 0+-----+ |
```

(continues on next page)

(continued from previous page)

```

| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03   CUDA Version: 12.
↵5      |
|-----+-----+-----+-----+-----+-----+-----+-----+-----+
↵-----+
| GPU   Name                     Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
↵ ECC |
| Fan  Temp   Perf           Pwr:Usage/Cap |      Memory-Usage | GPU-Util  ↵
↵Compute M. |
|
↵MIG M. |
|=====+=====+=====+=====+=====+=====+=====+=====+=====|
|   0   Tesla T4                       On   | 00000001:00:00.0 Off |                    ↵
↵ Off |
| N/A    42C    P8              9W / 70W |      2MiB / 16384MiB |           0%    ↵
↵Default |
|
↵ N/A |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
↵-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
↵-----+
| Processes:
↵
| GPU   GI    CI          PID    Type    Process name                        GPU
↵Memory |
|      ID    ID                                   Usage
↵
|=====+=====+=====+=====+=====+=====+=====+=====+=====|
| No running processes found
↵
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
↵-----+

```

## 6.2 Kesten processes

A **Kesten process** is a stochastic process of the form

$$X_{t+1} = a_{t+1}X_t + \eta_{t+1} \quad (6.1)$$

where  $\{a_t\}_{t \geq 1}$  and  $\{\eta_t\}_{t \geq 1}$  are IID sequences.

We are interested in the dynamics of  $\{X_t\}_{t \geq 0}$  when  $X_0$  is given.

We will focus on the nonnegative scalar case, where  $X_t$  takes values in  $\mathbb{R}_+$ .

In particular, we will assume that

- the initial condition  $X_0$  is nonnegative,
- $\{a_t\}_{t \geq 1}$  is a nonnegative IID stochastic process and
- $\{\eta_t\}_{t \geq 1}$  is another nonnegative IID stochastic process, independent of the first.

### 6.2.1 Application: firm dynamics

In this section we apply Kesten process theory to the study of firm dynamics.

#### Gibrat's law

It was postulated many years ago by Robert Gibrat that firm size evolves according to a simple rule whereby size next period is proportional to current size.

This is now known as [Gibrat's law of proportional growth](#).

We can express this idea by stating that a suitably defined measure  $s_t$  of firm size obeys

$$\frac{s_{t+1}}{s_t} = a_{t+1} \quad (6.2)$$

for some positive IID sequence  $\{a_t\}$ .

Subsequent empirical research has shown that this specification is not accurate, particularly for small firms.

However, we can get close to the data by modifying (6.2) to

$$s_{t+1} = a_{t+1}s_t + b_{t+1} \quad (6.3)$$

where  $\{a_t\}$  and  $\{b_t\}$  are both IID and independent of each other.

We now study the implications of this specification.

#### Heavy tails

If the conditions of the [Kesten–Goldie Theorem](#) are satisfied, then (6.3) implies that the firm size distribution will have Pareto tails.

This matches empirical findings across many data sets.

But there is another unrealistic aspect of the firm dynamics specified in (6.3) that we need to address: it ignores entry and exit.

In any given period and in any given market, we observe significant numbers of firms entering and exiting the market.

In this setting, firm dynamics can be expressed as

$$s_{t+1} = e_{t+1}\mathbb{1}\{s_t < \bar{s}\} + (a_{t+1}s_t + b_{t+1})\mathbb{1}\{s_t \geq \bar{s}\} \quad (6.4)$$

The motivation behind and interpretation of (6.2.4) can be found in [our earlier Kesten process lecture](#).

What can we say about dynamics?

Although (6.4) is not a Kesten process, it does update in the same way as a Kesten process when  $s_t$  is large.

So perhaps its stationary distribution still has Pareto tails?

We can investigate this question via simulation and rank-size plots.

The approach will be to

1. generate  $M$  draws of  $s_T$  when  $M$  and  $T$  are large and
2. plot the largest 1,000 of the resulting draws in a rank-size plot.

(The distribution of  $s_T$  will be close to the stationary distribution when  $T$  is large.)

In the simulation, we assume that each of  $a_t, b_t$  and  $e_t$  is lognormal.

Here's code to update a cross-section of firms according to the dynamics in (6.2.4).

```
@jax.jit
def update_s(s, s_bar, a_random, b_random, e_random):
    exp_a = jnp.exp(a_random)
    exp_b = jnp.exp(b_random)
    exp_e = jnp.exp(e_random)

    s = jnp.where(s < s_bar,
                  exp_e,
                  exp_a * s + exp_b)

    return s
```

Now we write a for loop that repeatedly calls this function, to push a cross-section of firms forward in time.

For sufficiently large  $T$ , the cross-section it returns (the cross-section at time  $T$ ) corresponds to firm size distribution in (approximate) equilibrium.

```
def generate_draws(M=1_000_000,
                  mu_a=-0.5,
                  sigma_a=0.1,
                  mu_b=0.0,
                  sigma_b=0.5,
                  mu_e=0.0,
                  sigma_e=0.5,
                  s_bar=1.0,
                  T=500,
                  s_init=1.0,
                  seed=123):

    key = random.PRNGKey(seed)

    # Initialize the array of s values with the initial value
    s = jnp.full((M, ), s_init)

    # Perform updates on s for time t
    for t in range(T):
        keys = random.split(key, 3)
        a_random = mu_a + sigma_a * random.normal(keys[0], (M, ))
        b_random = mu_b + sigma_b * random.normal(keys[1], (M, ))
        e_random = mu_e + sigma_e * random.normal(keys[2], (M, ))

        s = update_s(s, s_bar, a_random, b_random, e_random)

        # Generate new key for the next iteration
        key = random.fold_in(key, t)

    return s

%time data = generate_draws().block_until_ready()
```

```
2024-08-12 03:59:33.278172: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler, XLA will
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```

(continued from previous page)

```
CPU times: user 1.93 s, sys: 863 ms, total: 2.8 s
Wall time: 2.29 s
```

Running the above function again so we can see the speed with and without compile time.

```
%time data = generate_draws().block_until_ready()
```

```
CPU times: user 1.11 s, sys: 499 ms, total: 1.61 s
Wall time: 1.03 s
```

Notice that we do not JIT-compile the `for` loops, since

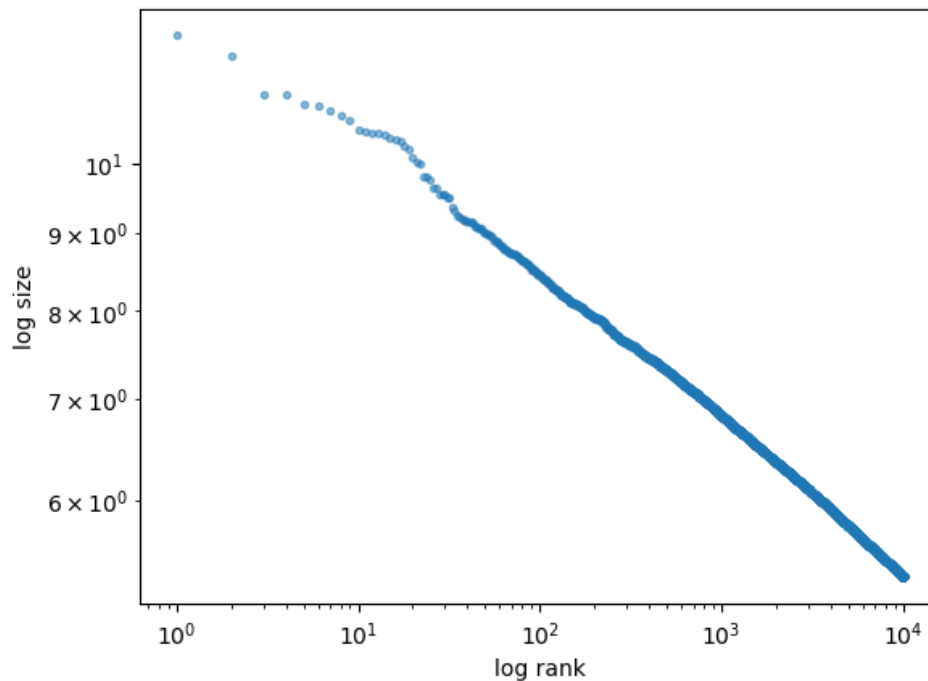
1. acceleration of the outer loop makes little difference terms of compute time and
2. compiling the outer loop is often very slow.

Let's produce the rank-size plot and check the distribution:

```
fig, ax = plt.subplots()

rank_data, size_data = qe.rank_size(data, c=0.01)
ax.loglog(rank_data, size_data, 'o', markersize=3.0, alpha=0.5)
ax.set_xlabel("log rank")
ax.set_ylabel("log size")

plt.show()
```



The plot produces a straight line, consistent with a Pareto tail.

### Alternative implementation with `lax.fori_loop`

If the time horizon is not too large, we can try to further accelerate our code by replacing the `for` loop with `lax.fori_loop`.

Note, however, that

1. as mentioned above, there is not much speed gain in accelerating outer loops,
2. `lax.fori_loop` has a more complicated syntax, and, most importantly,
3. the `lax.fori_loop` implementation consumes far more memory, as we need to have to store large matrices of random draws

Hence the code below will fail due to out-of-memory errors when  $T$  and  $M$  are large.

Here is the `lax.fori_loop` version:

```
@jax.jit
def generate_draws_lax(mu_a=-0.5,
                      sigma_a=0.1,
                      mu_b=0.0,
                      sigma_b=0.5,
                      mu_e=0.0,
                      sigma_e=0.5,
                      s_bar=1.0,
                      T=500,
                      M=500_000,
                      s_init=1.0,
                      seed=123):

    key = random.PRNGKey(seed)
    keys = random.split(key, 3)

    # Generate random draws and initial values
    a_random = mu_a + sigma_a * random.normal(keys[0], (T, M))
    b_random = mu_b + sigma_b * random.normal(keys[1], (T, M))
    e_random = mu_e + sigma_e * random.normal(keys[2], (T, M))
    s = jnp.full((M, ), s_init)

    # Define the function for each update
    def update_s(i, s):
        a, b, e = a_random[i], b_random[i], e_random[i]
        s = jnp.where(s < s_bar,
                      jnp.exp(e),
                      jnp.exp(a) * s + jnp.exp(b))

        return s

    # Use lax.scan to perform the calculations on all states
    s_final = lax.fori_loop(0, T, update_s, s)
    return s_final

%time data = generate_draws_lax().block_until_ready()
```

```
CPU times: user 584 ms, sys: 37.8 ms, total: 621 ms
Wall time: 629 ms
```

In this case,  $M$  is small enough for the code to run and we see some speed gain over the `for` loop implementation:

```
%time data = generate_draws_lax().block_until_ready()
```

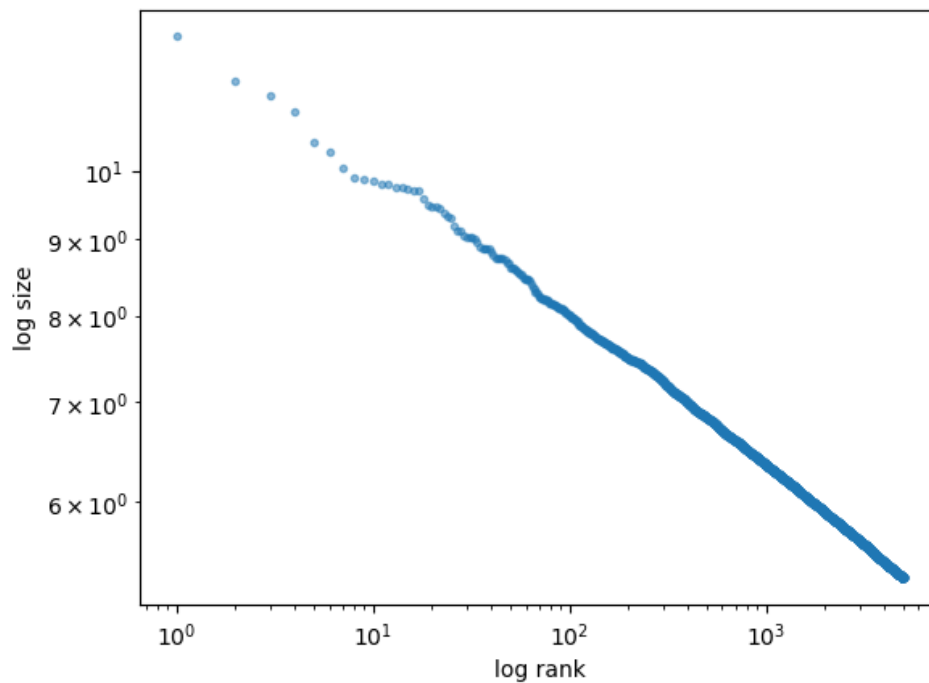
```
CPU times: user 3.31 ms, sys: 0 ns, total: 3.31 ms
Wall time: 60.3 ms
```

Here we produce the same rank-size plot:

```
fig, ax = plt.subplots()

rank_data, size_data = qe.rank_size(data, c=0.01)
ax.loglog(rank_data, size_data, 'o', markersize=3.0, alpha=0.5)
ax.set_xlabel("log rank")
ax.set_ylabel("log size")

plt.show()
```



Let's rerun the for loop version on smaller M to compare the speed

```
%time generate_draws(M=500_000).block_until_ready()
```

```
CPU times: user 1.52 s, sys: 564 ms, total: 2.09 s
Wall time: 1.53 s
```

```
Array([2.389801 , 2.2558599, 3.3113828, ..., 2.7102313, 2.5520844,
       3.4196172], dtype=float32)
```

Let's run it again to get rid of the compile time.

```
%time generate_draws(M=500_000).block_until_ready()
```

```
CPU times: user 1.19 s, sys: 363 ms, total: 1.56 s  
Wall time: 1e+03 ms
```

```
Array([2.389801 , 2.2558599, 3.3113828, ..., 2.7102313, 2.5520844,  
       3.4196172], dtype=float32)
```

We see that the `lax.fori_loop` version is faster than the `for` loop version when memory is not an issue.



## WEALTH DISTRIBUTION DYNAMICS

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

In this lecture we examine wealth dynamics in large cross-section of agents who are subject to both

- idiosyncratic shocks, which affect labor income and returns, and
- an aggregate shock, which also impacts on labor income and returns

In most macroeconomic models savings and consumption are determined by optimization.

Here savings and consumption behavior is taken as given – you can plug in your favorite model to obtain savings behavior and then analyze distribution dynamics using the techniques described below.

One of our interests will be how different aspects of wealth dynamics – such as labor income and the rate of return on investments – feed into measures of inequality, such as the Gini coefficient.

In addition to JAX and Anaconda, this lecture will need the following libraries:

```
!pip install quantecon
```

We will use the following imports:

```
import numba
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import quantecon as qe
import jax
import jax.numpy as jnp
from time import time
```

Let’s check the GPU we are running

```
!nvidia-smi
```

```

Mon Aug 12 04:05:58 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5 |
+-----+
+-----+
| GPU Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+=====+
|    0   Tesla T4                       On      | 00000001:00:00.0 Off  |
| Off |
| N/A    59C    P8              10W /  70W |      2MiB / 16384MiB |      0%
| Default |
|
| N/A |
+-----+
+-----+
+-----+
| Processes:
|
| GPU   GI    CI          PID    Type   Process name                      GPU
| Memory |
|      ID    ID
|
|
+=====+
| No running processes found
|
+-----+
+-----+

```

## 7.1 Wealth dynamics

Wealth evolves as follows:

$$w_{t+1} = (1 + r_{t+1})s(w_t) + y_{t+1}$$

Here

- $w_t$  is wealth at time  $t$  for a given household,
- $r_t$  is the rate of return of financial assets,
- $y_t$  is labor income and
- $s(w_t)$  is savings (current wealth minus current consumption)

There is an aggregate state process

$$z_{t+1} = az_t + b + \sigma_z \epsilon_{t+1}$$

that affects the interest rate and labor income.

In particular, the gross interest rates obey

$$R_t := 1 + r_t = c_r \exp(z_t) + \exp(\mu_r + \sigma_r \xi_t)$$

while

$$y_t = c_y \exp(z_t) + \exp(\mu_y + \sigma_y \zeta_t)$$

The tuple  $\{(\epsilon_t, \xi_t, \zeta_t)\}$  is IID and standard normal in  $\mathbb{R}^3$ .

(Each household receives their own idiosyncratic shocks.)

Regarding the savings function  $s$ , our default model will be

$$s(w) = s_0 w \cdot \mathbb{1}\{w \geq \hat{w}\} \quad (7.1)$$

where  $s_0$  is a positive constant.

Thus,

- for  $w < \hat{w}$ , the household saves nothing, while
- for  $w \geq \bar{w}$ , the household saves a fraction  $s_0$  of their wealth.

## 7.2 Implementation

### 7.2.1 Numba implementation

Here's a function that collects parameters and useful constants

```
def create_wealth_model(w_hat=1.0,      # Savings parameter
                        s_0=0.75,      # Savings parameter
                        c_y=1.0,        # Labor income parameter
                        mu_y=1.0,        # Labor income parameter
                        sigma_y=0.2,    # Labor income parameter
                        c_r=0.05,       # Rate of return parameter
                        mu_r=0.1,       # Rate of return parameter
                        sigma_r=0.5,    # Rate of return parameter
                        a=0.5,          # Aggregate shock parameter
                        b=0.0,          # Aggregate shock parameter
                        sigma_z=0.1):   # Aggregate shock parameter

    """
    Create a wealth model with given parameters.

    Return a tuple model = (household_params, aggregate_params), where
    household_params collects household information and aggregate_params
    collects information relevant to the aggregate shock process.

    """
    # Mean and variance of z process
    z_mean = b / (1 - a)
    z_var = sigma_z**2 / (1 - a**2)
    exp_z_mean = np.exp(z_mean + z_var / 2)
    # Mean of R and y processes
    R_mean = c_r * exp_z_mean + np.exp(mu_r + sigma_r**2 / 2)
```

(continues on next page)

(continued from previous page)

```

y_mean = c_y * exp_z_mean + np.exp(mu_y + sigma_y**2 / 2)
# Test stability condition ensuring wealth does not diverge
# to infinity.
alpha = R_mean * s_0
if alpha >= 1:
    raise ValueError("Stability condition failed.")
# Pack values into tuples and return them
household_params = (w_hat, s_0, c_y, mu_y, sigma_y, c_r, mu_r, sigma_r, y_mean)
aggregate_params = (a, b, sigma_z, z_mean, z_var)
model = household_params, aggregate_params
return model

```

Here's a function that generates the aggregate state process

```

@numba.jit
def generate_aggregate_state_sequence(aggregate_params, length=100):
    a, b, sigma_z, z_mean, z_var = aggregate_params
    z = np.empty(length+1)
    z[0] = z_mean # Initialize at z_mean
    for t in range(length):
        z[t+1] = a * z[t] + b + sigma_z * np.random.randn()
    return z

```

Here's a function that updates household wealth by one period, taking the current value of the aggregate shock

```

@numba.jit
def update_wealth(household_params, w, z):
    """
    Generate w_{t+1} given w_t and z_{t+1}.
    """
    # Unpack
    w_hat, s_0, c_y, mu_y, sigma_y, c_r, mu_r, sigma_r, y_mean = household_params
    # Update wealth
    y = c_y * np.exp(z) + np.exp(mu_y + sigma_y * np.random.randn())
    wp = y
    if w >= w_hat:
        R = c_r * np.exp(z) + np.exp(mu_r + sigma_r * np.random.randn())
        wp += R * s_0 * w
    return wp

```

Here's a function to simulate the time series of wealth for an individual household

```

@numba.jit
def wealth_time_series(model, w_0, sim_length):
    """
    Generate a single time series of length sim_length for wealth given initial
    value w_0. The function generates its own aggregate shock sequence.
    """
    # Unpack
    household_params, aggregate_params = model
    a, b, sigma_z, z_mean, z_var = aggregate_params
    # Initialize and update
    z = generate_aggregate_state_sequence(aggregate_params,
                                         length=sim_length)

```

(continues on next page)

(continued from previous page)

```

w = np.empty(sim_length)
w[0] = w_0
for t in range(sim_length-1):
    w[t+1] = update_wealth(household_params, w[t], z[t+1])
return w

```

Let's look at the wealth dynamics of an individual household

```

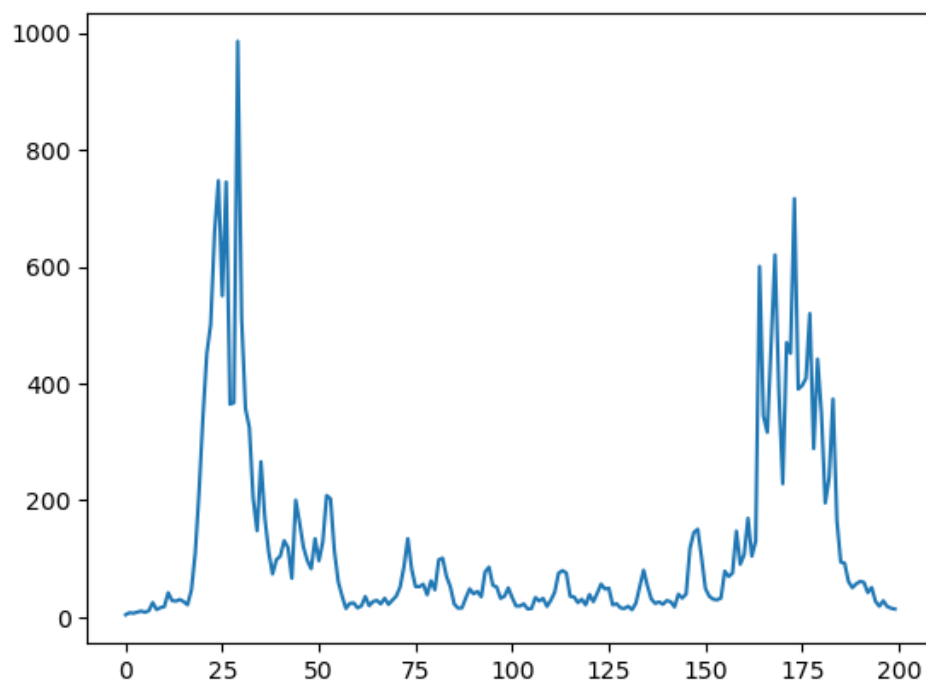
model = create_wealth_model()
household_params, aggregate_params = model
w_hat, s_0, c_y,  $\mu_y$ ,  $\sigma_y$ , c_r,  $\mu_r$ ,  $\sigma_r$ , y_mean = household_params
a, b,  $\sigma_z$ , z_mean, z_var = aggregate_params
ts_length = 200
w = wealth_time_series(model, y_mean, ts_length)

```

```

fig, ax = plt.subplots()
ax.plot(w)
plt.show()

```



Notice the large spikes in wealth over time.

Such spikes are related to heavy tails in the wealth distribution, which we discuss below.

Here's a function to simulate a cross section of households forward in time.

Note the use of parallelization to speed up computation.

```

@numba.jit(parallel=True)
def update_cross_section(model, w_distribution, z_sequence):
    """
    Shifts a cross-section of households forward in time
    """

```

(continues on next page)

(continued from previous page)

```

Takes

    * a current distribution of wealth values as w_distribution and
    * an aggregate shock sequence z_sequence

and updates each w_t in w_distribution to w_{t+j}, where
j = len(z_sequence).

Returns the new distribution.

"""
# Unpack
household_params, aggregate_params = model

num_households = len(w_distribution)
new_distribution = np.empty_like(w_distribution)
z = z_sequence

# Update each household
for i in numba.prange(num_households):
    w = w_distribution[i]
    for t in range(sim_length):
        w = update_wealth(household_params, w, z[t])
    new_distribution[i] = w
return new_distribution

```

Parallelization works in the function above because the time path of each household can be calculated independently once the path for the aggregate state is known.

Let's see how long it takes to shift a large cross-section of households forward 200 periods

```

sim_length = 200
num_households = 10_000_000
psi_0 = np.full(num_households, y_mean) # Initial distribution
z_sequence = generate_aggregate_state_sequence(aggregate_params,
                                              length=sim_length)

print("Generating cross-section using Numba")
start = time()
psi_star = update_cross_section(model, psi_0, z_sequence)
numba_with_compile = time() - start
print(f"Generated cross-section in {numba_with_compile} seconds.\n")

```

```
Generating cross-section using Numba
```

```
Generated cross-section in 57.59461736679077 seconds.
```

We run it again to eliminate compile time.

```

start = time()
psi_star = update_cross_section(model, psi_0, z_sequence)
numba_without_compile = time() - start
print(f"Generated cross-section in {numba_without_compile} seconds.\n")

```

Generated cross-section in 56.80775547027588 seconds.

## 7.2.2 JAX implementation

Let's redo some of the preceding calculations using JAX and see how execution speed compares

```
def update_cross_section_jax(model, w_distribution, z_sequence, key):
    """
    Shifts a cross-section of households forward in time

    Takes

        * a current distribution of wealth values as w_distribution and
        * an aggregate shock sequence z_sequence

    and updates each w_t in w_distribution to w_{t+j}, where
    j = len(z_sequence).

    Returns the new distribution.

    """
    # Unpack, simplify names
    household_params, aggregate_params = model
    w_hat, s_0, c_y, mu_y, sigma_y, c_r, mu_r, sigma_r, y_mean = household_params
    w = w_distribution
    n = len(w)

    # Update wealth
    for t, z in enumerate(z_sequence):
        U = jax.random.normal(key, (2, n))
        y = c_y * jnp.exp(z) + jnp.exp(mu_y + sigma_y * U[0, :])
        R = c_r * jnp.exp(z) + jnp.exp(mu_r + sigma_r * U[1, :])
        w = y + jnp.where(w < w_hat, 0.0, R * s_0 * w)
        key, subkey = jax.random.split(key)

    return w
```

Let's see how long it takes to shift the cross-section of households forward using JAX

```
sim_length = 200
num_households = 10_000_000
psi_0 = jnp.full(num_households, y_mean) # Initial distribution
z_sequence = generate_aggregate_state_sequence(aggregate_params,
                                                length=sim_length)
z_sequence = jnp.array(z_sequence)
```

```
2024-08-12 04:07:54.613303: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```

```
print("Generating cross-section using JAX")
key = jax.random.PRNGKey(1234)
start = time()
ψ_star = update_cross_section_jax(model, ψ_0, z_sequence, key).block_until_ready()
jax_with_compile = time() - start
print(f"Generated cross-section in {jax_with_compile} seconds.\n")
```

Generating cross-section using JAX

Generated cross-section in 2.3496782779693604 seconds.

```
print("Repeating without compile time.")
key = jax.random.PRNGKey(1234)
start = time()
ψ_star = update_cross_section_jax(model, ψ_0, z_sequence, key).block_until_ready()
jax_without_compile = time() - start
print(f"Generated cross-section in {jax_without_compile} seconds")
```

Repeating without compile time.

Generated cross-section in 1.235994577407837 seconds

And let's see how long it takes if we compile the loop.

```
def update_cross_section_jax_compiled(model,
                                      w_distribution,
                                      w_size,
                                      z_sequence,
                                      key):
    """
    Shifts a cross-section of households forward in time

    Takes

        * a current distribution of wealth values as w_distribution and
        * an aggregate shock sequence z_sequence

    and updates each w_t in w_distribution to w_{t+j}, where
    j = len(z_sequence).

    Returns the new distribution.

    """
    # Unpack, simplify names
    household_params, aggregate_params = model
    w_hat, s_0, c_y, μ_y, σ_y, c_r, μ_r, σ_r, y_mean = household_params
    w = w_distribution
    n = len(w)
    z = z_sequence
    sim_length = len(z)

    def body_function(t, state):
        key, w = state
```

(continues on next page)



(continued from previous page)

```

key, subkey = jax.random.split(key)
U = jax.random.normal(subkey, (2, n))
y = c_y * jnp.exp(z[t]) + jnp.exp(mu_y + sigma_y * U[0, :])
R = c_r * jnp.exp(z[t]) + jnp.exp(mu_r + sigma_r * U[1, :])
w = y + jnp.where(w < w_hat, 0.0, R * s_0 * w)
return key, w

key, w = jax.lax.fori_loop(0, sim_length, body_function, (key, w))
return w

```

```

update_cross_section_jax_compiled = jax.jit(
    update_cross_section_jax_compiled, static_argnums=(2,)
)

```

```

print("Generating cross-section using JAX with compiled loop")
key = jax.random.PRNGKey(1234)
start = time()
psi_star = update_cross_section_jax_compiled(
    model, psi_0, num_households, z_sequence, key
).block_until_ready()
jax_for_i_with_compile = time() - start
print(f"Generated cross-section in {jax_for_i_with_compile} seconds.\n")

```

Generating cross-section using JAX with compiled loop

Generated cross-section in 0.570561408996582 seconds.

```

print("Repeating without compile time")
key = jax.random.PRNGKey(1234)
start = time()
psi_star = update_cross_section_jax_compiled(
    model, psi_0, num_households, z_sequence, key
).block_until_ready()
jax_for_i_without_compile = time() - start
print(f"Generated cross-section in {jax_for_i_without_compile} seconds")

```

Repeating without compile time  
Generated cross-section in 0.1854405403137207 seconds

```

print(f"JAX is {numba_without_compile/jax_for_i_without_compile:.4f} times faster.\n")

```

JAX is 306.3395 times faster.

### 7.2.3 Pareto tails

In most countries, the cross-sectional distribution of wealth exhibits a Pareto tail (power law).

Let's see if our model can replicate this stylized fact by running a simulation that generates a cross-section of wealth and generating a suitable rank-size plot.

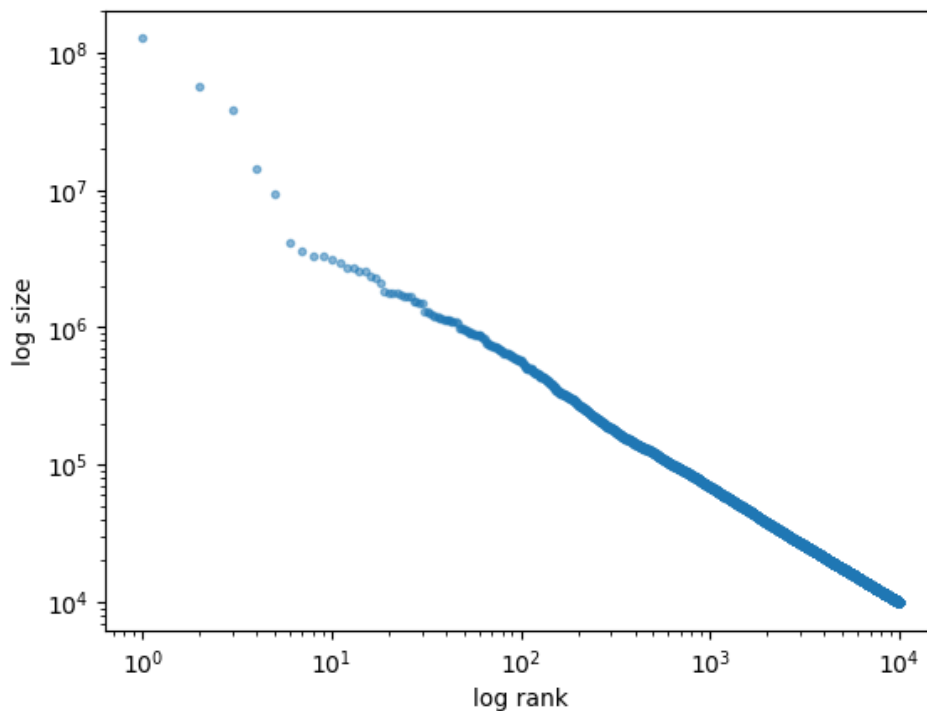
We will use the function `rank_size` from `quantecon` library.

In the limit, data that obeys a power law generates a straight line.

```
model = create_wealth_model()
key = jax.random.PRNGKey(1234)
ψ_star = update_cross_section_jax_compiled(
    model, ψ_0, num_households, z_sequence, key
)
fig, ax = plt.subplots()

rank_data, size_data = qe.rank_size(ψ_star, c=0.001)
ax.loglog(rank_data, size_data, 'o', markersize=3.0, alpha=0.5)
ax.set_xlabel("log rank")
ax.set_ylabel("log size")

plt.show()
```



## 7.2.4 Lorenz curves and Gini coefficients

To study the impact of parameters on inequality, we examine Lorenz curves and the Gini coefficients at different parameters.

QuantEcon provides functions to compute Lorenz curves and Gini coefficients that are accelerated using Numba.

Here we provide JAX-based functions that do the same job and are faster for large data sets on parallel hardware.

### Lorenz curve

Recall that, for sorted data  $w_1, \dots, w_n$ , the Lorenz curve generates data points  $(x_i, y_i)_{i=0}^n$  according to

$$x_0 = y_0 = 0 \quad \text{and, for } i \geq 1, \quad x_i = \frac{i}{n}, \quad y_i = \frac{\sum_{j \leq i} w_j}{\sum_{j \leq n} w_j}$$

```
def _lorenz_curve_jax(w, w_size):
    n = w.shape[0]
    w = jnp.sort(w)
    x = jnp.arange(n + 1) / n
    s = jnp.concatenate((jnp.zeros(1), jnp.cumsum(w)))
    y = s / s[n]
    return x, y

lorenz_curve_jax = jax.jit(_lorenz_curve_jax, static_argnums=(1,))
```

Let's test

```
sim_length = 200
num_households = 1_000_000
psi_0 = jnp.full(num_households, y_mean) # Initial distribution
z_sequence = generate_aggregate_state_sequence(aggregate_params,
                                                length=sim_length)
z_sequence = jnp.array(z_sequence)
```

```
key = jax.random.PRNGKey(1234)
psi_star = update_cross_section_jax_compiled(
    model, psi_0, num_households, z_sequence, key
)
```

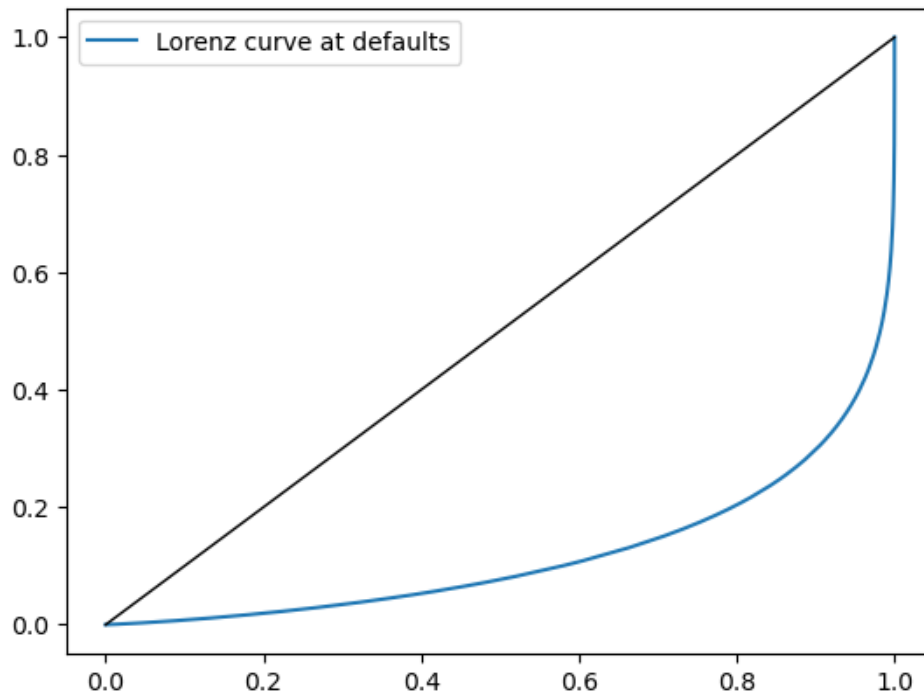
```
%time _ = lorenz_curve_jax(psi_star, num_households)
```

```
CPU times: user 936 ms, sys: 4.02 ms, total: 940 ms
Wall time: 1.39 s
```

```
# Now time it without compile time
%time x, y = lorenz_curve_jax(psi_star, num_households)
```

```
CPU times: user 378 µs, sys: 1 µs, total: 379 µs
Wall time: 197 µs
```

```
fig, ax = plt.subplots()
ax.plot(x, y, label="Lorenz curve at defaults")
ax.plot(x, x, 'k-', lw=1)
ax.legend()
plt.show()
```



### Gini Coefficient

Recall that, for sorted data  $w_1, \dots, w_n$ , the Gini coefficient takes the form

$$G := \frac{\sum_{i=1}^n \sum_{j=1}^n |w_j - w_i|}{2n \sum_{i=1}^n w_i}. \quad (7.2)$$

Here's a function that computes the Gini coefficient using vectorization.

```
def _gini_jax(w, w_size):
    w_1 = jnp.reshape(w, (w_size, 1))
    w_2 = jnp.reshape(w, (1, w_size))
    g_sum = jnp.sum(jnp.abs(w_1 - w_2))
    return g_sum / (2 * w_size * jnp.sum(w))

gini_jax = jax.jit(_gini_jax, static_argnums=(1,))
```

```
%time gini = gini_jax(ψ_star, num_households).block_until_ready()
```

```
CPU times: user 79.3 ms, sys: 9 μs, total: 79.3 ms
Wall time: 8.16 s
```

```
# Now time it without compilation
%time gini = gini_jax( $\psi$ _star, num_households).block_until_ready()
```

```
CPU times: user 1.37 ms, sys: 1 ms, total: 2.37 ms
Wall time: 7.99 s
```

```
gini
```

```
Array(0.7594847, dtype=float32)
```

## 7.3 Exercises

### Exercise 7.3.1

In this exercise, write an alternative version of `gini_jax` that uses `vmap` instead of reshaping and broadcasting. Test with the same array to see if you can obtain the same output

### Solution to Exercise 7.3.1

Here's one solution:

```
@jax.jit
def gini_jax_vmap(w):

    def _inner_sum(x):
        return jnp.sum(jnp.abs(x - w))

    inner_sum = jax.vmap(_inner_sum)

    full_sum = jnp.sum(inner_sum(w))
    return full_sum / (2 * len(w) * jnp.sum(w))
```

```
%time gini = gini_jax_vmap( $\psi$ _star).block_until_ready()
```

```
CPU times: user 79.6 ms, sys: 11  $\mu$ s, total: 79.7 ms
Wall time: 8.16 s
```

```
# Now time it without compile time
%time gini = gini_jax_vmap( $\psi$ _star).block_until_ready()
```

```
CPU times: user 2.41 ms, sys: 7  $\mu$ s, total: 2.42 ms
Wall time: 8.09 s
```

```
gini
```

```
Array(0.7594847, dtype=float32)
```

---

### Exercise 7.3.2

In this exercise we investigate how the parameters determining the rate of return on assets and labor income shape inequality.

In doing so we recall that

$$R_t := 1 + r_t = c_r \exp(z_t) + \exp(\mu_r + \sigma_r \xi_t)$$

while

$$y_t = c_y \exp(z_t) + \exp(\mu_y + \sigma_y \zeta_t)$$

Investigate how the Lorenz curves and the Gini coefficient associated with the wealth distribution change as return to savings varies.

In particular, plot Lorenz curves for the following three different values of  $\mu_r$ .

```
 $\mu_r$ _vals = (0.0, 0.025, 0.05)
```

Use the following as your initial cross-sectional distribution

```
num_households = 1_000_000
 $\psi_0$  = jnp.full(num_households, y_mean) # Initial distribution
```

Once you have done that, plot the Gini coefficients as well.

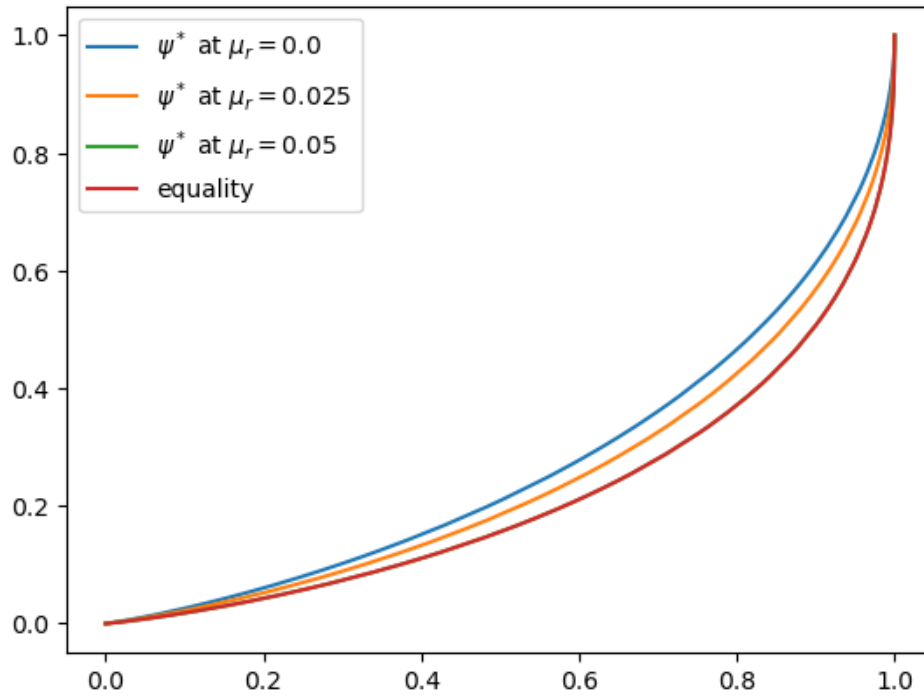
Do the outcomes match your intuition?

---

### Solution to Exercise 7.3.2

Here is one solution

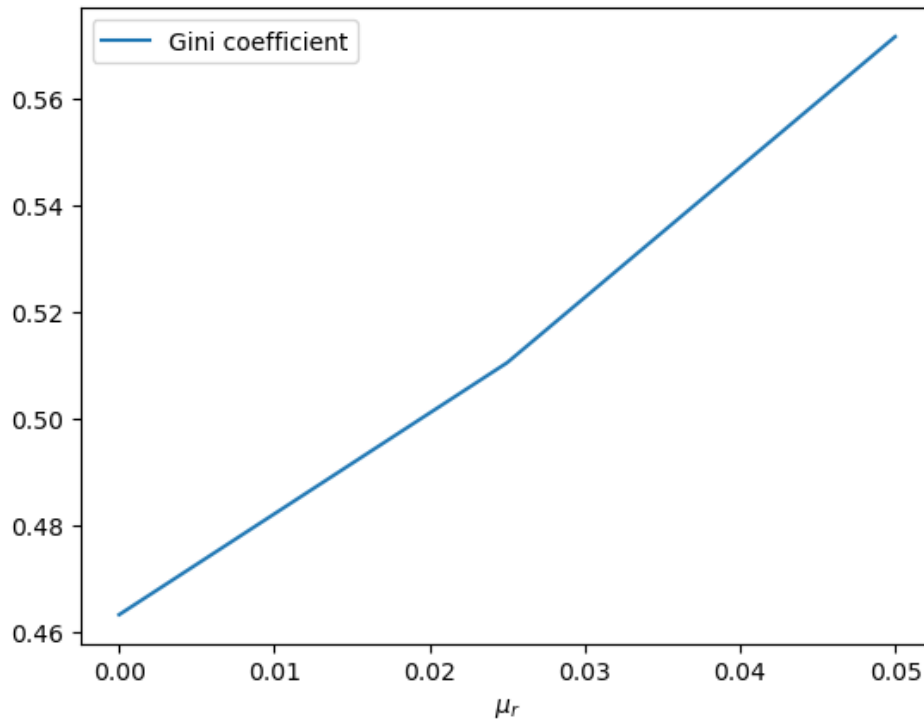
```
key = jax.random.PRNGKey(1234)
fig, ax = plt.subplots()
gini_vals = []
for  $\mu_r$  in  $\mu_r$ _vals:
    model = create_wealth_model( $\mu_r$ = $\mu_r$ )
     $\psi_{\text{star}}$  = update_cross_section_jax_compiled(
        model,  $\psi_0$ , num_households, z_sequence, key
    )
    x, y = lorenz_curve_jax( $\psi_{\text{star}}$ , num_households)
    g = gini_jax( $\psi_{\text{star}}$ , num_households)
    ax.plot(x, y, label=f'$\psi^{\wedge}$ at $\mu_r = \{\mu_r:0.2\}$')
    gini_vals.append(g)
ax.plot(x, y, label='equality')
ax.legend(loc="upper left")
plt.show()
```



The Lorenz curve shifts downwards as returns on financial income rise, indicating a rise in inequality.

Now let's check the Gini coefficient

```
fig, ax = plt.subplots()
ax.plot(p_r_vals, gini_vals, label='Gini coefficient')
ax.set_xlabel(" $\mu_r$ ")
ax.legend()
plt.show()
```



As expected, inequality increases as returns on financial income rise.

---

### Exercise 7.3.3

Now investigate what happens when we change the volatility term  $\sigma_r$  in financial returns.

Use the same initial condition as before and the sequence

```
σ_r_vals = (0.35, 0.45, 0.52)
```

To isolate the role of volatility, set  $\mu_r = -\sigma_r^2/2$  at each  $\sigma_r$ .

(This holds the variance of the idiosyncratic term  $\exp(\mu_r + \sigma_r \zeta)$  constant.)

---

### Solution to Exercise 7.3.3

Here's one solution

```
key = jax.random.PRNGKey(1234)
fig, ax = plt.subplots()

gini_vals = []
for σ_r in σ_r_vals:
    model = create_wealth_model(σ_r=σ_r, μ_r=(-σ_r**2/2))
    ψ_star = update_cross_section_jax_compiled(
        model, ψ_0, num_households, z_sequence, key
    )
    x, y = lorenz_curve_jax(ψ_star, num_households)
    g = gini_jax(ψ_star, num_households)
```

(continues on next page)

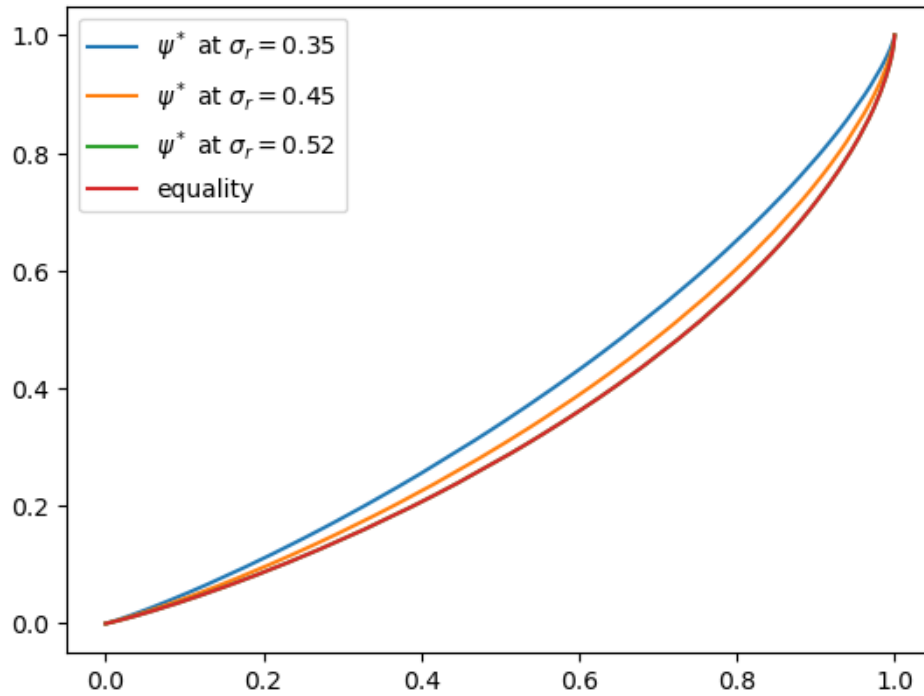


(continued from previous page)

```

ax.plot(x, y, label=f'$\psi^*$ at $\sigma_r = \{\sigma_r:0.2\}$')
gini_vals.append(g)
ax.plot(x, y, label='equality')
ax.legend(loc="upper left")
plt.show()

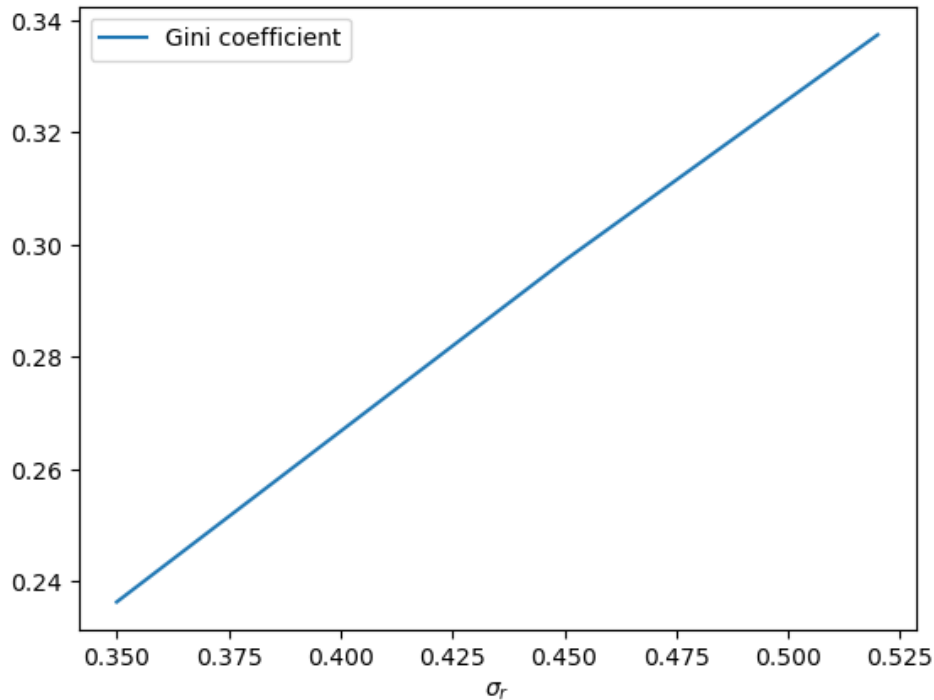
```



```

fig, ax = plt.subplots()
ax.plot(sigma_r_vals, gini_vals, label='Gini coefficient')
ax.set_xlabel("$\sigma_r$")
ax.legend()
plt.show()

```



---

**Exercise 7.3.4**

In this exercise, examine which has more impact on inequality:

- a 5% rise in volatility of the rate of return,
- or a 5% rise in volatility of labor income.

Test this by

1. Shifting  $\sigma_r$  up 5% from the baseline and plotting the Lorenz curve
2. Shifting  $\sigma_y$  up 5% from the baseline and plotting the Lorenz curve

Plot both on the same figure and examine the result.

---

**Solution to Exercise 7.3.4**

Here's one solution.

It shows that increasing volatility in financial income has a greater effect

```
model = create_wealth_model()
household_params, aggregate_params = model
w_hat, s_0, c_y, mu_y, sigma_y, c_r, mu_r, sigma_r, y_mean = household_params
sigma_r_default = sigma_r
sigma_y_default = sigma_y

psi_star = update_cross_section_jax_compiled(
    model, psi_0, num_households, z_sequence, key
)
```

(continues on next page)

(continued from previous page)

```

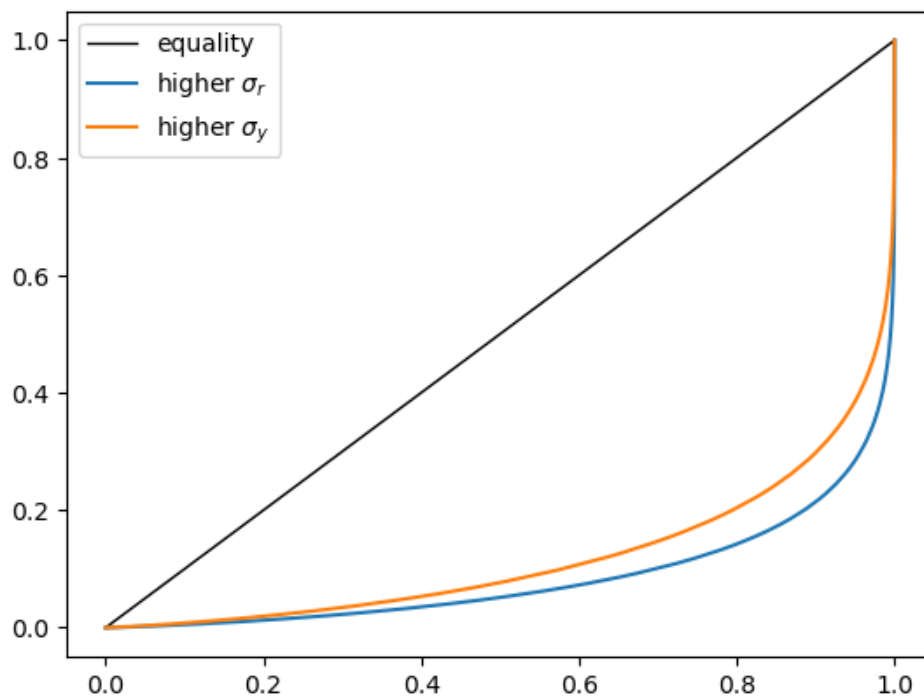
x_default, y_default = lorenz_curve_jax( $\psi$ _star, num_households)

model = create_wealth_model( $\sigma_r$ =(1.05 *  $\sigma_r$ _default))
 $\psi$ _star = update_cross_section_jax_compiled(
    model,  $\psi$ _0, num_households, z_sequence, key
)
x_financial, y_financial = lorenz_curve_jax( $\psi$ _star, num_households)

model = create_wealth_model( $\sigma_y$ =(1.05 *  $\sigma_y$ _default))
 $\psi$ _star = update_cross_section_jax_compiled(
    model,  $\psi$ _0, num_households, z_sequence, key
)
x_labor, y_labor = lorenz_curve_jax( $\psi$ _star, num_households)

fig, ax = plt.subplots()
ax.plot(x_default, x_default, 'k-', lw=1, label='equality')
ax.plot(x_financial, y_financial, label=r'higher  $\sigma_r$ ')
ax.plot(x_labor, y_labor, label=r'higher  $\sigma_y$ ')
ax.legend()
plt.show()

```





# **Part III**

## **Asset Pricing**



## ASSET PRICING: THE LUCAS ASSET PRICING MODEL

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

### 8.1 Overview

An asset is a claim on a stream of prospective payments.

What is the correct price to pay for such a claim?

The asset pricing model of Lucas [[Lucas, 1978](#)] attempts to answer this question in an equilibrium setting with risk-averse agents.

Lucas’ model provides a beautiful illustration of model building in general and equilibrium pricing in competitive models in particular.

In this lecture we work through the Lucas model and show where the fundamental asset pricing equation comes from.

We’ll write code in both Numba and JAX.

Since the model is relatively small, the speed gain from using JAX is not as large as it is in some other lectures.

Nonetheless, the gain is nontrivial.

Let’s start with some imports:

```
import jax.numpy as jnp
import jax
import numpy as np
import numba
from scipy.stats import lognorm
import matplotlib.pyplot as plt
from time import time
```

## 8.2 The Lucas Model

Lucas studied a pure exchange economy with a representative consumer (or household), where

- *Pure exchange* means that all endowments are exogenous.
- *Representative* consumer means that either
  - there is a single consumer (sometimes also referred to as a household), or
  - all consumers have identical endowments and preferences

Either way, the assumption of a representative agent means that prices adjust to eradicate desires to trade.

This makes it very easy to compute competitive equilibrium prices.

### 8.2.1 Basic Setup

Let's review the setup.

#### Assets

There is a single “productive unit” that costlessly generates a sequence of consumption goods  $\{y_t\}_{t=0}^{\infty}$ .

Another way to view  $\{y_t\}_{t=0}^{\infty}$  is as a *consumption endowment* for this economy.

We will assume that this endowment is Markovian, following the exogenous process

$$y_{t+1} = G(y_t, \xi_{t+1})$$

Here  $\{\xi_t\}$  is an IID shock sequence with known distribution  $\phi$  and  $y_t \geq 0$ .

An asset is a claim on all or part of this endowment stream.

The consumption goods  $\{y_t\}_{t=0}^{\infty}$  are nonstorable, so holding assets is the only way to transfer wealth into the future.

For the purposes of intuition, it's common to think of the productive unit as a “tree” that produces fruit.

Based on this idea, a “Lucas tree” is a claim on the consumption endowment.

#### Consumers

A representative consumer ranks consumption streams  $\{c_t\}$  according to the time separable utility functional

$$\mathbb{E} \sum_{t=0}^{\infty} \beta^t u(c_t) \tag{8.1}$$

Here

- $\beta \in (0, 1)$  is a fixed discount factor.
- $u$  is a strictly increasing, strictly concave, continuously differentiable period utility function.
- $\mathbb{E}$  is a mathematical expectation.



## 8.2.2 Pricing a Lucas Tree

What is an appropriate price for a claim on the consumption endowment?

We'll price an *ex-dividend* claim, meaning that

- the seller retains this period's dividend
- the buyer pays  $p_t$  today to purchase a claim on
  - $y_{t+1}$  and
  - the right to sell the claim tomorrow at price  $p_{t+1}$

Since this is a competitive model, the first step is to pin down consumer behavior, taking prices as given.

Next, we'll impose equilibrium constraints and try to back out prices.

In the consumer problem, the consumer's control variable is the share  $\pi_t$  of the claim held in each period.

Thus, the consumer problem is to maximize (8.1) subject to

$$c_t + \pi_{t+1}p_t \leq \pi_t y_t + \pi_t p_t$$

along with  $c_t \geq 0$  and  $0 \leq \pi_t \leq 1$  at each  $t$ .

The decision to hold share  $\pi_t$  is actually made at time  $t - 1$ .

But this value is inherited as a state variable at time  $t$ , which explains the choice of subscript.

### The Dynamic Program

We can write the consumer problem as a dynamic programming problem.

Our first observation is that prices depend on current information, and current information is really just the endowment process up until the current period.

In fact, the endowment process is Markovian, so that the only relevant information is the current state  $y \in \mathbb{R}_+$  (dropping the time subscript).

This leads us to guess an equilibrium where price is a function  $p$  of  $y$ .

Remarks on the solution method

- Since this is a competitive (read: price taking) model, the consumer will take this function  $p$  as given.
- In this way, we determine consumer behavior given  $p$  and then use equilibrium conditions to recover  $p$ .
- This is the standard way to solve competitive equilibrium models.

Using the assumption that price is a given function  $p$  of  $y$ , we write the value function and constraint as

$$v(\pi, y) = \max_{c, \pi'} \left\{ u(c) + \beta \int v(\pi', G(y, z)) \phi(dz) \right\}$$

subject to

$$c + \pi' p(y) \leq \pi y + \pi p(y) \tag{8.2}$$

We can invoke the fact that utility is increasing to claim equality in (8.2) and hence eliminate the constraint, obtaining

$$v(\pi, y) = \max_{\pi'} \left\{ u[\pi(y + p(y)) - \pi' p(y)] + \beta \int v(\pi', G(y, z)) \phi(dz) \right\} \tag{8.3}$$

The solution to this dynamic programming problem is an optimal policy expressing either  $\pi'$  or  $c$  as a function of the state  $(\pi, y)$ .

- Each one determines the other, since  $c(\pi, y) = \pi(y + p(y)) - \pi'(\pi, y)p(y)$

## Next Steps

What we need to do now is determine equilibrium prices.

It seems that to obtain these, we will have to

1. Solve this two-dimensional dynamic programming problem for the optimal policy.
2. Impose equilibrium constraints.
3. Solve out for the price function  $p(y)$  directly.

However, as Lucas showed, there is a related but more straightforward way to do this.

## Equilibrium Constraints

Since the consumption good is not storable, in equilibrium we must have  $c_t = y_t$  for all  $t$ .

In addition, since there is one representative consumer (alternatively, since all consumers are identical), there should be no trade in equilibrium.

In particular, the representative consumer owns the whole tree in every period, so  $\pi_t = 1$  for all  $t$ .

Prices must adjust to satisfy these two constraints.

## The Equilibrium Price Function

Now observe that the first-order condition for (8.3) can be written as

$$u'(c)p(y) = \beta \int v'_1(\pi', G(y, z))\phi(dz)$$

where  $v'_1$  is the derivative of  $v$  with respect to its first argument.

To obtain  $v'_1$  we can simply differentiate the right-hand side of (8.3) with respect to  $\pi$ , yielding

$$v'_1(\pi, y) = u'(c)(y + p(y))$$

Next, we impose the equilibrium constraints while combining the last two equations to get

$$p(y) = \beta \int \frac{u'[G(y, z)]}{u'(y)} [G(y, z) + p(G(y, z))] \phi(dz) \quad (8.4)$$

In sequential rather than functional notation, we can also write this as

$$p_t = \mathbb{E}_t \left[ \beta \frac{u'(c_{t+1})}{u'(c_t)} (y_{t+1} + p_{t+1}) \right] \quad (8.5)$$

This is the famous consumption-based asset pricing equation.

Before discussing it further we want to solve out for prices.

### 8.2.3 Solving the Model

Equation (8.4) is a *functional equation* in the unknown function  $p$ .

The solution is an equilibrium price function  $p^*$ .

Let's look at how to obtain it.

#### Setting up the Problem

Instead of solving for it directly we'll follow Lucas' indirect approach, first setting

$$f(y) := u'(y)p(y) \quad (8.6)$$

so that (8.4) becomes

$$f(y) = h(y) + \beta \int f[G(y, z)]\phi(dz) \quad (8.7)$$

Here  $h(y) := \beta \int u'[G(y, z)]G(y, z)\phi(dz)$  is a function that depends only on the primitives.

Equation (8.7) is a functional equation in  $f$ .

The plan is to solve out for  $f$  and convert back to  $p$  via (8.6).

To solve (8.7) we'll use a standard method: convert it to a fixed point problem.

First, we introduce the operator  $T$  mapping  $f$  into  $Tf$  as defined by

$$(Tf)(y) = h(y) + \beta \int f[G(y, z)]\phi(dz) \quad (8.8)$$

In what follows, we refer to  $T$  as the Lucas operator.

The reason we do this is that a solution to (8.7) now corresponds to a function  $f^*$  satisfying  $(Tf^*)(y) = f^*(y)$  for all  $y$ .

In other words, a solution is a *fixed point* of  $T$ .

This means that we can use fixed point theory to obtain and compute the solution.

#### A Little Fixed Point Theory

Let  $cb\mathbb{R}_+$  be the set of continuous bounded functions  $f: \mathbb{R}_+ \rightarrow \mathbb{R}_+$ .

We now show that

1.  $T$  has exactly one fixed point  $f^*$  in  $cb\mathbb{R}_+$ .
2. For any  $f \in cb\mathbb{R}_+$ , the sequence  $T^k f$  converges uniformly to  $f^*$ .

---

**Note:** If you find the mathematics heavy going you can take 1–2 as given and skip to the [next section](#)

---

Recall the [Banach contraction mapping theorem](#).

It tells us that the previous statements will be true if we can find an  $\alpha < 1$  such that

$$\|Tf - Tg\| \leq \alpha \|f - g\|, \quad \forall f, g \in cb\mathbb{R}_+ \quad (8.9)$$

Here  $\|h\| := \sup_{x \in \mathbb{R}_+} |h(x)|$ .

To see that (8.9) is valid, pick any  $f, g \in cb\mathbb{R}_+$  and any  $y \in \mathbb{R}_+$ .

Observe that, since integrals get larger when absolute values are moved to the inside,

$$\begin{aligned} |Tf(y) - Tg(y)| &= \left| \beta \int f[G(y, z)] \phi(dz) - \beta \int g[G(y, z)] \phi(dz) \right| \\ &\leq \beta \int |f[G(y, z)] - g[G(y, z)]| \phi(dz) \\ &\leq \beta \int \|f - g\| \phi(dz) \\ &= \beta \|f - g\| \end{aligned}$$

Since the right-hand side is an upper bound, taking the sup over all  $y$  on the left-hand side gives (8.9) with  $\alpha := \beta$ .

## 8.3 Computation

The preceding discussion tells that we can compute  $f^*$  by picking any arbitrary  $f \in cb\mathbb{R}_+$  and then iterating with  $T$ .

The equilibrium price function  $p^*$  can then be recovered by  $p^*(y) = f^*(y)/u'(y)$ .

Let's try this when  $\ln y_{t+1} = \alpha \ln y_t + \sigma \epsilon_{t+1}$  where  $\{\epsilon_t\}$  is IID and standard normal.

Utility will take the isoelastic form  $u(c) = c^{1-\gamma}/(1-\gamma)$ , where  $\gamma > 0$  is the coefficient of relative risk aversion.

We'll use Monte Carlo to compute the integral

$$\int f[G(y, z)] \phi(dz)$$

Monte Carlo is not always the fastest method for computing low-dimensional integrals, but it is extremely flexible (for example, it's straightforward to change the underlying state process).

### 8.3.1 Numba Code

Let's start with code using NumPy / Numba (and then compare it to code using JAX).

We create a function that returns tuples containing parameters and arrays needed for computation.

```
def create_lucas_tree_model( $\gamma=2$ ,          # CRRA utility parameter
                            $\beta=0.95$ ,      # Discount factor
                            $\alpha=0.90$ ,     # Correlation coefficient
                            $\sigma=0.1$ ,      # Volatility coefficient
                           grid_size=500,
                           draw_size=1_000,
                           seed=11):
    # Set the grid interval to contain most of the mass of the
    # stationary distribution of the consumption endowment
     $ssd = \sigma / \text{np.sqrt}(1 - \alpha^2)$ 
    grid_min, grid_max =  $\text{np.exp}(-4 * ssd)$ ,  $\text{np.exp}(4 * ssd)$ 
    grid =  $\text{np.linspace}(grid\_min, grid\_max, grid\_size)$ 
    # Set up distribution for shocks
     $\text{np.random.seed}(seed)$ 
     $\phi = \text{lognorm}(\sigma)$ 
    draws =  $\phi.\text{rvs}(500)$ 
    # And the vector h
```

(continues on next page)

(continued from previous page)

```

h = np.empty(grid_size)
for i, y in enumerate(grid):
    h[i] =  $\beta$  * np.mean((y** $\alpha$  * draws)**(1 -  $\gamma$ ))
# Pack and return
params =  $\gamma$ ,  $\beta$ ,  $\alpha$ ,  $\sigma$ 
arrays = grid, draws, h
return params, arrays

```

Here's a Numba-jitted version of the Lucas operator

```

@numba.jit
def T(params, arrays, f):
    """
    The Lucas pricing operator.
    """
    # Unpack
     $\gamma$ ,  $\beta$ ,  $\alpha$ ,  $\sigma$  = params
    grid, draws, h = arrays
    # Turn f into a function
    Af = lambda x: np.interp(x, grid, f)
    # Compute Tf and return
    Tf = np.empty_like(f)
    # Apply the T operator to f using Monte Carlo integration
    for i in range(len(grid)):
        y = grid[i]
        Tf[i] = h[i] +  $\beta$  * np.mean(Af(y** $\alpha$  * draws))
    return Tf

```

To solve the model, we write a function that iterates using the Lucas operator to find the fixed point.

```

def solve_model(params, arrays, tol=1e-6, max_iter=500):
    """
    Compute the equilibrium price function.
    """
    # Unpack
     $\gamma$ ,  $\beta$ ,  $\alpha$ ,  $\sigma$  = params
    grid, draws, h = arrays
    # Set up and loop
    i = 0
    f = np.ones_like(grid) # Initial guess of f
    error = tol + 1
    while error > tol and i < max_iter:
        Tf = T(params, arrays, f)
        error = np.max(np.abs(Tf - f))
        f = Tf
        i += 1
    price = f * grid** $\gamma$  # Back out price vector
    return price

```

Let's solve the model and plot the resulting price function

```

params, arrays = create_lucas_tree_model()
 $\gamma$ ,  $\beta$ ,  $\alpha$ ,  $\sigma$  = params
grid, draws, h = arrays

```

(continues on next page)

(continued from previous page)

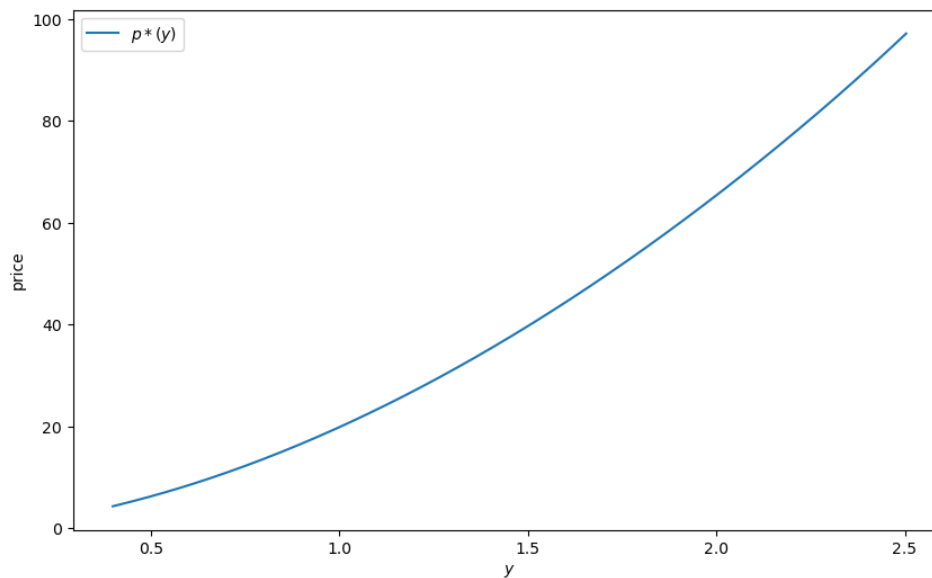
```
# Solve once to compile
start = time()
price_vals = solve_model(params, arrays)
numba_with_compile_time = time() - start
print("Numba compile plus execution time = ", numba_with_compile_time)
```

```
Numba compile plus execution time = 4.966493129730225
```

```
# Now time execution without compile time
start = time()
price_vals = solve_model(params, arrays)
numba_without_compile_time = time() - start
print("Numba execution time = ", numba_without_compile_time)
```

```
Numba execution time = 3.333735942840576
```

```
fig, ax = plt.subplots(figsize=(10, 6))
ax.plot(grid, price_vals, label='$p^*(y)$')
ax.set_xlabel('$y$')
ax.set_ylabel('price')
ax.legend()
plt.show()
```



We see that the price is increasing, even if we remove all serial correlation from the endowment process.

The reason is that a larger current endowment reduces current marginal utility.

The price must therefore rise to induce the household to consume the entire endowment (and hence satisfy the resource constraint).

### 8.3.2 JAX Code

Here's a JAX version of the same problem.

```
def create_lucas_tree_model(y=2,          # CRRA utility parameter
                           β=0.95,      # Discount factor
                           α=0.90,      # Correlation coefficient
                           σ=0.1,        # Volatility coefficient
                           grid_size=500,
                           draw_size=1_000,
                           seed=11):
    # Set the grid interval to contain most of the mass of the
    # stationary distribution of the consumption endowment
    ssd = σ / jnp.sqrt(1 - α**2)
    grid_min, grid_max = jnp.exp(-4 * ssd), jnp.exp(4 * ssd)
    grid = jnp.linspace(grid_min, grid_max, grid_size)

    # Set up distribution for shocks
    key = jax.random.key(seed)
    draws = jax.random.lognormal(key, σ, shape=(draw_size,))
    grid_resaped = grid.reshape((grid_size, 1))
    draws_resaped = draws.reshape((-1, draw_size))
    h = β * jnp.mean((grid_resaped**α * draws_resaped) ** (1-y), axis=1)
    params = y, β, α, σ
    arrays = grid, draws, h
    return params, arrays
```

We'll use the following function to simultaneously compute the expectation

$$\int f[G(y, z)]\phi(dz)$$

over all  $y$  in the grid, under the current specifications.

```
@jax.jit
def compute_expectation(y, α, draws, grid, f):
    return jnp.mean(jnp.interp(y**α * draws, grid, f))

# Vectorize over y
compute_expectation = jax.vmap(compute_expectation,
                               in_axes=(0, None, None, None, None))
```

Here's the Lucas operator

```
@jax.jit
def T(params, arrays, f):
    """
    The Lucas operator

    """
    grid, draws, h = arrays
    y, β, α, σ = params
    mci = compute_expectation(grid, α, draws, grid, f)
    return h + β * mci
```

We'll use successive approximation to compute the fixed point.

```

def successive_approx_jax(T,                # Operator (callable)
                        x_0,                # Initial condition
                        tol=1e-6,           # Error tolerance
                        max_iter=10_000):    # Max iteration bound
    def body_fun(k_x_err):
        k, x, error = k_x_err
        x_new = T(x)
        error = jnp.max(jnp.abs(x_new - x))
        return k + 1, x_new, error

    def cond_fun(k_x_err):
        k, x, error = k_x_err
        return jnp.logical_and(error > tol, k < max_iter)

    k, x, error = jax.lax.while_loop(cond_fun, body_fun,
                                     (1, x_0, tol + 1))

    return x

successive_approx_jax = \
    jax.jit(successive_approx_jax, static_argnums=(0,))

```

Here's a function that solves the model

```

def solve_model(params, arrays, tol=1e-6, max_iter=500):
    """
    Compute the equilibrium price function.

    """
    # Simplify notation
    grid, draws, h = arrays
    y, β, α, σ = params
    _T = lambda f: T(params, arrays, f)
    f = jnp.ones_like(grid) # Initial guess of f

    f = successive_approx_jax(_T, f, tol=tol, max_iter=max_iter)

    price = f * grid**y # Back out price vector

    return price

```

Now let's solve the model again and compare timing

```

params, arrays = create_lucas_tree_model()
grid, draws, h = arrays
y, β, α, σ = params

# Solve once to compile
start = time()
price_vals = solve_model(params, arrays).block_until_ready()
jax_with_compile_time = time() - start
print("JAX compile plus execution time = ", jax_with_compile_time)

```

```

2024-08-12 03:59:51.987426: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA for Linux
↳compatibility packages.
(continues on next page)

```



(continued from previous page)

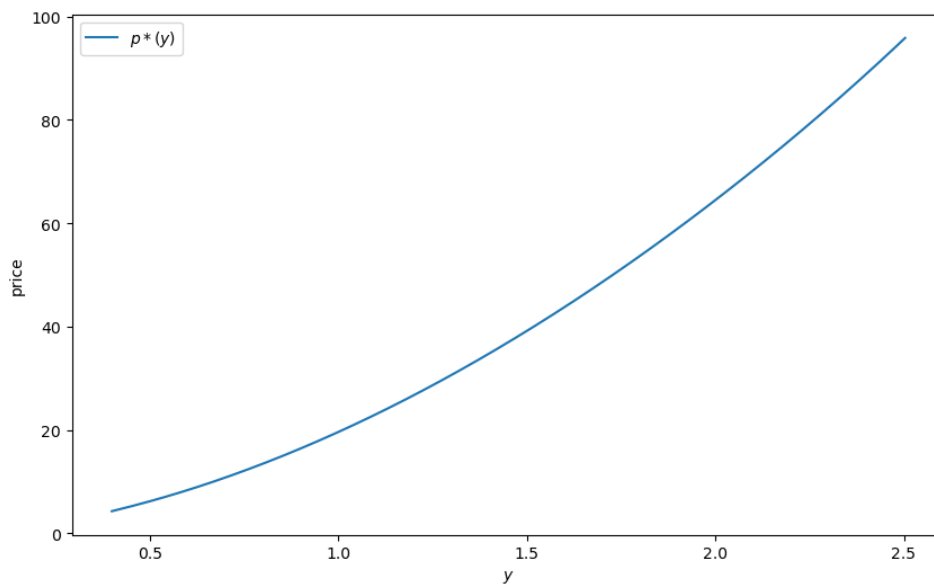
```
JAX compile plus execution time = 0.7722537517547607
```

```
# Now time execution without compile time
start = time()
price_vals = solve_model(params, arrays).block_until_ready()
jax_without_compile_time = time() - start
print("JAX execution time = ", jax_without_compile_time)
print("Speedup factor = ", numba_without_compile_time/jax_without_compile_time)
```

```
JAX execution time = 0.4955613613128662
Speedup factor = 6.72719102637194
```

Let's check the solutions are similar

```
fig, ax = plt.subplots(figsize=(10, 6))
ax.plot(grid, price_vals, label='$p^*(y)$')
ax.set_xlabel('$y$')
ax.set_ylabel('price')
ax.legend()
plt.show()
```



## 8.4 Exercises

---

### Exercise 8.4.1

When consumers are more patient the asset becomes more valuable, and the price of the Lucas tree shifts up.

Show this by plotting the price function for the Lucas tree when  $\beta = 0.95$  and  $0.98$ .

---

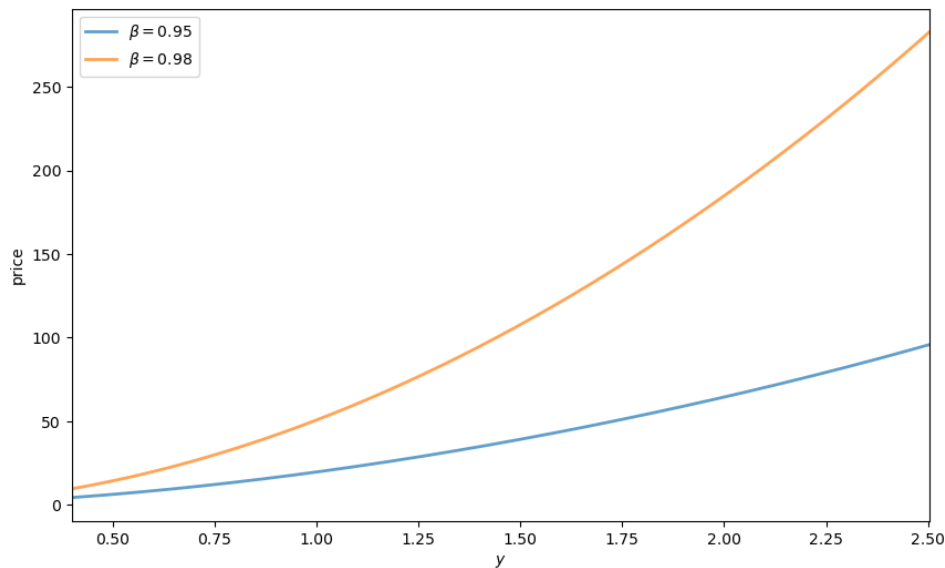
---

### Solution to Exercise 8.4.1

```
fig, ax = plt.subplots(figsize=(10, 6))

for  $\beta$  in (.95, 0.98):
    params, arrays = create_lucas_tree_model( $\beta$ = $\beta$ )
    grid, draws, h = arrays
     $\gamma$ , beta,  $\alpha$ ,  $\sigma$  = params
    price_vals = solve_model(params, arrays)
    label = rf'$\beta = {beta}$'
    ax.plot(grid, price_vals, lw=2, alpha=0.7, label=label)

ax.legend(loc='upper left')
ax.set(xlabel='$y$', ylabel='price', xlim=(min(grid), max(grid)))
plt.show()
```



## AN ASSET PRICING PROBLEM

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

### 9.1 Overview

In this lecture we consider some asset pricing problems and use them to illustrate some foundations of JAX programming.

The main difference from the lecture *Asset Pricing: The Lucas Asset Pricing Model*, which also considers asset prices, is that the state spaces will be discrete and multi-dimensional.

Most of the heavy lifting is done through routines from linear algebra.

Along the way, we will show how to solve some memory-intensive problems with large state spaces.

We do this using elegant techniques made available by JAX, involving the use of linear operators to avoid instantiating large matrices.

If you wish to skip all motivation and move straight to the first equation we plan to solve, you can jump to (9.5.5).

The code outputs below are generated by machine connected to the following GPU

```
!nvidia-smi
```

```
Mon Aug 12 03:59:57 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5                                    |
+-----+-----+-----+-----+
| GPU   Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC   |
| Fan  Temp  Perf    Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
+-----+-----+-----+-----+
|MIG M. |
```

(continues on next page)



Roughly speaking, (9.2.1) says that the cost (i.e., price) equals expected benefit.

The discount factor is introduced because most people prefer payments now to payments in the future.

One problem with this very simple model is that it does not take into account attitudes to risk.

For example, investors often demand higher rates of return for holding risky assets.

This feature of asset prices cannot be captured by risk neutral pricing.

Hence we modify (9.2.1) to

$$P_t = \mathbb{E}_t M_{t+1} D_{t+1} \quad (9.2)$$

In this expression,  $M_{t+1}$  replaces  $\beta$  and is called the **stochastic discount factor**.

In essence, allowing discounting to become a random variable gives us the flexibility to combine temporal discounting and attitudes to risk.

We leave further discussion to [other lectures](#) because our aim is to move to the computational problem.

## 9.3 Pricing a cash flow

Now let's try to price an asset like a share, which delivers a cash flow  $D_t, D_{t+1}, \dots$

We will call these payoffs “dividends”.

If we buy the share, hold it for one period and sell it again, we receive one dividend and our payoff is  $D_{t+1} + P_{t+1}$ .

Therefore, by (9.2.2), the price should be

$$P_t = \mathbb{E}_t M_{t+1} [D_{t+1} + P_{t+1}] \quad (9.3)$$

Because prices generally grow over time, which complicates analysis, it will be easier for us to solve for the **price-dividend ratio**  $V_t := P_t/D_t$ .

Let's write down an expression that this ratio should satisfy.

We can divide both sides of (9.3) by  $D_t$  to get

$$V_t = \mathbb{E}_t \left[ M_{t+1} \frac{D_{t+1}}{D_t} (1 + V_{t+1}) \right] \quad (9.4)$$

We can also write this as

$$V_t = \mathbb{E}_t [M_{t+1} \exp(G_{t+1}^d) (1 + V_{t+1})] \quad (9.5)$$

where

$$G_{t+1}^d = \ln \frac{D_{t+1}}{D_t}$$

is the growth rate of dividends.

Our aim is to solve (9.3.3) but before that we need to specify

1. the stochastic discount factor  $M_{t+1}$  and
2. the growth rate of dividends  $G_{t+1}^d$

## 9.4 Choosing the stochastic discount factor

We will adopt the stochastic discount factor described in *Asset Pricing: The Lucas Asset Pricing Model*, which has the form

$$M_{t+1} = \beta \frac{u'(C_{t+1})}{u'(C_t)} \quad (9.6)$$

where  $u$  is a utility function and  $C_t$  is time  $t$  consumption of a representative consumer.

For utility, we'll assume the **constant relative risk aversion** (CRRA) specification

$$u(c) = \frac{c^{1-\gamma}}{1-\gamma} \quad (9.7)$$

Inserting the CRRA specification into (9.6) and letting

$$G_{t+1}^c = \ln \left( \frac{C_{t+1}}{C_t} \right)$$

the growth rate rate of consumption, we obtain

$$M_{t+1} = \beta \left( \frac{C_{t+1}}{C_t} \right)^{-\gamma} = \beta \exp(G_{t+1}^c)^{-\gamma} = \beta \exp(-\gamma G_{t+1}^c) \quad (9.8)$$

## 9.5 Solving for the price-dividend ratio

Substituting (9.4.3) into (9.5) gives the price-dividend ratio formula

$$V_t = \beta \mathbb{E}_t [\exp(G_{t+1}^d - \gamma G_{t+1}^c)(1 + V_{t+1})] \quad (9.9)$$

We assume there is a Markov chain  $\{X_t\}$ , which we call the **state process**, such that

$$\begin{aligned} G_{t+1}^c &= \mu_c + X_t + \sigma_c \epsilon_{c,t+1} \\ G_{t+1}^d &= \mu_d + X_t + \sigma_d \epsilon_{d,t+1} \end{aligned}$$

Here  $\{\epsilon_{c,t}\}$  and  $\{\epsilon_{d,t}\}$  are IID and standard normal, and independent of each other.

We can think of  $\{X_t\}$  as an aggregate shock that affects both consumption growth and firm profits (and hence dividends).

We let  $P$  be the **stochastic matrix that governs  $\{X_t\}$**  and assume  $\{X_t\}$  takes values in some finite set  $S$ .

We guess that  $V_t$  is a fixed function of this state process (and this guess turns out to be correct).

This means that  $V_t = v(X_t)$  for some unknown function  $v$ .

By (9.5.1), the unknown function  $v$  satisfies the equation

$$v(X_t) = \beta \mathbb{E}_t \left\{ \exp[a + (1-\gamma)X_t + \sigma_d \epsilon_{d,t+1} - \gamma \sigma_c \epsilon_{c,t+1}] (1 + v(X_{t+1})) \right\} \quad (9.10)$$

where  $a := \mu_d - \gamma \mu_c$

Since the shocks  $\epsilon_{c,t+1}$  and  $\epsilon_{d,t+1}$  are independent of  $\{X_t\}$ , we can integrate them out.

We use the following property of lognormal distributions: if  $Y = \exp(c\epsilon)$  for constant  $c$  and  $\epsilon \sim N(0, 1)$ , then  $\mathbb{E}Y = \exp(c^2/2)$ .

This yields

$$v(X_t) = \beta \mathbb{E}_t \left\{ \exp \left[ a + (1 - \gamma)X_t + \frac{\sigma_d^2 + \gamma^2 \sigma_c^2}{2} \right] (1 + v(X_{t+1})) \right\} \quad (9.11)$$

Conditioning on  $X_t = x$ , we can write this as

$$v(x) = \beta \sum_{y \in S} \left\{ \exp \left[ a + (1 - \gamma)x + \frac{\sigma_d^2 + \gamma^2 \sigma_c^2}{2} \right] (1 + v(y)) \right\} P(x, y) \quad (9.12)$$

for all  $x \in S$ .

Suppose  $S = \{x_1, \dots, x_N\}$ .

Then we can think of  $v$  as an  $N$ -vector and, using square brackets for indices on arrays, write

$$v[i] = \beta \sum_{j=1}^N \left\{ \exp \left[ a + (1 - \gamma)x[i] + \frac{\sigma_d^2 + \gamma^2 \sigma_c^2}{2} \right] (1 + v[j]) \right\} P[i, j] \quad (9.13)$$

for  $i = 1, \dots, N$ .

Equivalently, we can write

$$v[i] = \sum_{j=1}^N K[i, j](1 + v[j]) \quad (9.14)$$

where  $K$  is the matrix defined by

$$K[i, j] = \beta \left\{ \exp \left[ a + (1 - \gamma)x[i] + \frac{\sigma_d^2 + \gamma^2 \sigma_c^2}{2} \right] \right\} P[i, j] \quad (9.15)$$

Rewriting (9.5.6) in vector form yields

$$v = K(\mathbf{1} + v) \quad (9.16)$$

Notice that (9.5.8) can be written as  $(I - K)v = K\mathbf{1}$ .

The Neumann series lemma tells us that  $I - K$  is invertible and the solution is

$$v = (I - K)^{-1} K\mathbf{1} \quad (9.17)$$

whenever  $r(K)$ , the spectral radius of  $K$ , is strictly less than one.

Once we specify  $P$  and all the parameters, we can

1. obtain  $K$
2. check the spectral radius condition  $r(K) < 1$  and, assuming it holds,
3. compute the solution via (9.5.9).

## 9.6 Code

We will use the [power iteration algorithm](#) to check the spectral radius condition.

The function below computes the spectral radius of  $A$ .

```
def power_iteration_sr(A, num_iterations=15, seed=1234):
    """ Estimates the spectral radius of A via power iteration. """

    # Initialize
    key = jax.random.PRNGKey(seed)
    b_k = jax.random.normal(key, (A.shape[1],))
    sr = 0

    for _ in range(num_iterations):
        # calculate the matrix-by-vector product Ab
        b_k1 = jnp.dot(A, b_k)

        # calculate the norm
        b_k1_norm = jnp.linalg.norm(b_k1)

        # Record the current estimate of the spectral radius
        sr = jnp.sum(b_k1 * b_k) / jnp.sum(b_k * b_k)

        # re-normalize the vector and continue
        b_k = b_k1 / b_k1_norm

    return sr

power_iteration_sr = jax.jit(power_iteration_sr)
```

The next function verifies that the spectral radius of a given matrix is  $< 1$ .

```
def test_stability(Q):
    """
    Assert that the spectral radius of matrix Q is < 1.
    """
    sr = power_iteration_sr(Q)
    assert sr < 1, f"Spectral radius condition failed with radius = {sr}"
```

In what follows we assume that  $\{X_t\}$ , the state process, is a discretization of the AR(1) process

$$X_{t+1} = \rho X_t + \sigma \eta_{t+1}$$

where  $\rho, \sigma$  are parameters and  $\{\eta_t\}$  is IID and standard normal.

To discretize this process we use `QuantEcon.py`'s `tauchen` function.

Below we write a function called `create_model()` that returns a namedtuple storing the relevant parameters and arrays.

```
Model = namedtuple('Model',
                  ('P', 'S', 'β', 'γ', 'μ_c', 'μ_d', 'σ_c', 'σ_d'))

def create_model(N=100,                # size of state space for Markov chain
                 ρ=0.9,                 # persistence parameter for Markov chain
                 σ=0.01,                 # persistence parameter for Markov chain
                 β=0.98,                 # discount factor
                 γ=2.5,                  # coefficient of risk aversion
                 μ_c=0.01,               # mean growth of consumption
                 μ_d=0.01,               # mean growth of dividends
                 σ_c=0.02,               # consumption volatility
                 σ_d=0.04):              # dividend volatility
```

(continues on next page)



(continued from previous page)

```

# Create the state process
mc = qe.tauchen(N, ρ, σ)
S = mc.state_values
P = mc.P
# Shift arrays to the device
S, P = map(jax.device_put, (S, P))
# Return the namedtuple
return Model(P=P, S=S, β=β, γ=γ, μ_c=μ_c, μ_d=μ_d, σ_c=σ_c, σ_d=σ_d)

```

Our first step is to construct the matrix  $K$  defined in (9.5.7).

Here's a function that does this using loops.

```

def compute_K_loop(model):
    # unpack
    P, S, β, γ, μ_c, μ_d, σ_c, σ_d = model
    N = len(S)
    K = np.empty((N, N))
    a = μ_d - γ * μ_c
    for i, x in enumerate(S):
        for j, y in enumerate(S):
            e = np.exp(a + (1 - γ) * x + (σ_d**2 + γ**2 * σ_c**2) / 2)
            K[i, j] = β * e * P[i, j]
    return K

```

To exploit the parallelization capabilities of JAX, let's also write a vectorized (i.e., loop-free) implementation.

```

def compute_K(model):
    # unpack
    P, S, β, γ, μ_c, μ_d, σ_c, σ_d = model
    N = len(S)
    # Reshape and multiply pointwise using broadcasting
    x = np.reshape(S, (N, 1))
    a = μ_d - γ * μ_c
    e = np.exp(a + (1 - γ) * x + (σ_d**2 + γ**2 * σ_c**2) / 2)
    K = β * e * P
    return K

```

These two functions produce the same output:

```

model = create_model(N=10)
K1 = compute_K(model)
K2 = compute_K_loop(model)
np.allclose(K1, K2)

```

```

2024-08-12 04:00:02.351552: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.

```

```
True
```

Now we can compute the price-dividend ratio:

```

def price_dividend_ratio(model, test_stable=True):
    """
    Computes the price-dividend ratio of the asset.

    Parameters
    -----
    model: an instance of Model
            contains primitives

    Returns
    -----
    v : array_like
        price-dividend ratio

    """
    K = compute_K(model)
    N = len(model.S)

    if test_stable:
        test_stability(K)

    # Compute v
    I = np.identity(N)
    ones_vec = np.ones(N)
    v = np.linalg.solve(I - K, K @ ones_vec)

    return v

```

Here's a plot of  $v$  as a function of the state for several values of  $\gamma$ .

```

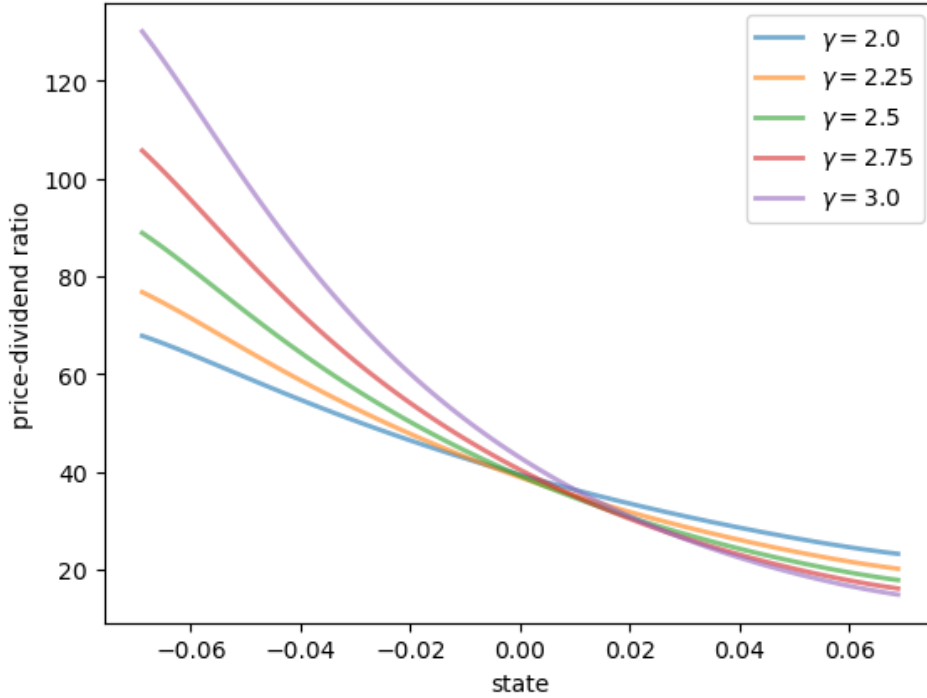
model = create_model()
S = model.S
ys = np.linspace(2.0, 3.0, 5)

fig, ax = plt.subplots()

for y in ys:
    model = create_model(y=y)
    v = price_dividend_ratio(model)
    ax.plot(S, v, lw=2, alpha=0.6, label=rf"$\gamma = {y}$")

ax.set_ylabel("price-dividend ratio")
ax.set_xlabel("state")
ax.legend(loc='upper right')
plt.show()

```



Notice that  $v$  is decreasing in each case.

This is because, with a positively correlated state process, higher states indicate higher future consumption growth.

With the stochastic discount factor (9.8), higher growth decreases the discount factor, lowering the weight placed on future dividends.

## 9.7 An Extended Example

One problem with the last set is that volatility is constant through time (i.e.,  $\sigma_c$  and  $\sigma_d$  are constants).

In reality, financial markets and growth rates of macroeconomic variables exhibit bursts of volatility.

To accommodate this, we now develop a *stochastic volatility* model.

To begin, suppose that consumption and dividends grow as follows.

$$G_{t+1}^i = \mu_i + Z_t + \bar{\sigma} \exp(H_t^i \epsilon_{i,t+1}), \quad i \in \{c, d\}$$

where  $\{Z_t\}$  is a finite Markov chain and  $\{H_t^c\}$  and  $\{H_t^d\}$  are volatility processes.

We assume that  $\{H_t^c\}$  and  $\{H_t^d\}$  are AR(1) processes of the form

$$H_{t+1}^i = \rho_i H_t^i + \sigma_i \eta_{i,t+1}, \quad i \in \{c, d\}$$

Here  $\{\eta_t^c\}$  and  $\{\eta_t^d\}$  are IID and standard normal.

Let  $X_t = (H_t^c, H_t^d, Z_t)$ .

We call  $\{X_t\}$  the state process and guess that  $V_t$  is a function of this state process, so that  $V_t = v(X_t)$  for some unknown function  $v$ .

Modifying (9.5.2) to accommodate the new growth specifications, we find that  $v$  satisfies

$$v(X_t) = \beta \times \mathbb{E}_t \left\{ \exp[a + (1 - \gamma)Z_t + \bar{\sigma} \exp(H_t^d)\epsilon_{d,t+1} - \gamma \bar{\sigma} \exp(H_t^c)\epsilon_{c,t+1}] (1 + v(X_{t+1})) \right\} \quad (9.18)$$

where, as before,  $a := \mu_d - \gamma\mu_c$

Conditioning on state  $x = (h_c, h_d, z)$ , this becomes

$$v(x) = \beta \mathbb{E}_t \exp[a + (1 - \gamma)z + \bar{\sigma} \exp(h_d)\epsilon_{d,t+1} - \gamma \bar{\sigma} \exp(h_c)\epsilon_{c,t+1}] (1 + v(X_{t+1})) \quad (9.19)$$

As before, we integrate out the independent shocks and use the rules for expectations of lognormals to obtain

$$v(x) = \beta \mathbb{E}_t \exp \left[ a + (1 - \gamma)z + \bar{\sigma}^2 \frac{\exp(2h_d) + \gamma^2 \exp(2h_c)}{2} \right] (1 + v(X_{t+1})) \quad (9.20)$$

Let

$$A(h_c, h_d, z, h'_c, h'_d, z') := \beta \exp \left[ a + (1 - \gamma)z + \bar{\sigma}^2 \frac{\exp(2h_d) + \gamma^2 \exp(2h_c)}{2} \right] P(h_c, h'_c) Q(h_d, h'_d) R(z, z')$$

where  $P, Q, R$  are the stochastic matrices for, respectively, discretized  $\{H_t^c\}$ , discretized  $\{H_t^d\}$  and  $\{Z_t\}$ ,

With this notation, we can write (9.7.3) more explicitly as

$$v(h_c, h_d, z) = \sum_{h'_c, h'_d, z'} (1 + v(h'_c, h'_d, z')) A(h_c, h_d, z, h'_c, h'_d, z') \quad (9.21)$$

Let's now write the state using indices, with  $(i, j, k)$  being the indices for  $(h_c, h_d, z)$ .

Then (9.21) becomes

$$v[i, j, k] = \sum_{i', j', k'} A[i, j, k, i', j', k'] (1 + v[i', j', k']) \quad (9.22)$$

One way to understand this is to reshape  $v$  into an  $N$ -vector, where  $N = I \times J \times K$ , and  $A$  into an  $N \times N$  matrix.

Then we can write (9.7.5) as

$$v = A(\mathbb{1} + v)$$

Provided that the spectral radius condition  $r(A) < 1$  holds, the solution is given by

$$v = (I - A)^{-1} A \mathbb{1}$$

## 9.8 Numpy Version

Our first implementation will be in NumPy.

Once we have a NumPy version working, we will convert it to JAX and check the difference in the run times.

The code block below provides a function called `create_sv_model()` that returns a namedtuple containing arrays and other data that form the primitives of the problem.

It assumes that  $\{Z_t\}$  is a discretization of

$$Z_{t+1} = \rho_z Z_t + \sigma_z \xi_{t+1}$$

```

SVModel = namedtuple('SVModel',
                    ('P', 'hc_grid',
                     'Q', 'hd_grid',
                     'R', 'z_grid',
                     'β', 'γ', 'bar_σ', 'μ_c', 'μ_d'))

def create_sv_model(β=0.98,          # discount factor
                   γ=2.5,           # coefficient of risk aversion
                   I=14,            # size of state space for h_c
                   ρ_c=0.9,         # persistence parameter for h_c
                   σ_c=0.01,        # volatility parameter for h_c
                   J=14,            # size of state space for h_d
                   ρ_d=0.9,         # persistence parameter for h_d
                   σ_d=0.01,        # volatility parameter for h_d
                   K=14,            # size of state space for z
                   bar_σ=0.01,      # volatility scaling parameter
                   ρ_z=0.9,         # persistence parameter for z
                   σ_z=0.01,        # persistence parameter for z
                   μ_c=0.001,       # mean growth of consumption
                   μ_d=0.005):      # mean growth of dividends

    mc = qe.tauchen(I, ρ_c, σ_c)
    hc_grid = mc.state_values
    P = mc.P
    mc = qe.tauchen(J, ρ_d, σ_d)
    hd_grid = mc.state_values
    Q = mc.P
    mc = qe.tauchen(K, ρ_z, σ_z)
    z_grid = mc.state_values
    R = mc.P

    return SVModel(P=P, hc_grid=hc_grid,
                   Q=Q, hd_grid=hd_grid,
                   R=R, z_grid=z_grid,
                   β=β, γ=γ, bar_σ=bar_σ, μ_c=μ_c, μ_d=μ_d)

```

Now we provide a function to compute the matrix  $A$ .

```

def compute_A(sv_model):
    # Set up
    P, hc_grid, Q, hd_grid, R, z_grid, β, γ, bar_σ, μ_c, μ_d = sv_model
    I, J, K = len(hc_grid), len(hd_grid), len(z_grid)
    N = I * J * K
    # Reshape and broadcast over (i, j, k, i', j', k')
    hc = np.reshape(hc_grid, (I, 1, 1, 1, 1, 1))
    hd = np.reshape(hd_grid, (1, J, 1, 1, 1, 1))
    z = np.reshape(z_grid, (1, 1, K, 1, 1, 1))
    P = np.reshape(P, (I, 1, 1, I, 1, 1))
    Q = np.reshape(Q, (1, J, 1, 1, J, 1))
    R = np.reshape(R, (1, 1, K, 1, 1, K))
    # Compute A and then reshape to create a matrix
    a = μ_d - γ * μ_c
    b = bar_σ**2 * (np.exp(2 * hd) + γ**2 * np.exp(2 * hc)) / 2
    κ = np.exp(a + (1 - γ) * z + b)
    A = β * κ * P * Q * R
    A = np.reshape(A, (N, N))
    return A

```

Here's our function to compute the price-dividend ratio for the stochastic volatility model.

```
def sv_pd_ratio(sv_model, test_stable=True):
    """
    Computes the price-dividend ratio of the asset for the stochastic volatility
    model.

    Parameters
    -----
    sv_model: an instance of Model
              contains primitives

    Returns
    -----
    v : array_like
        price-dividend ratio

    """
    # unpack
    P, hc_grid, Q, hd_grid, R, z_grid,  $\beta$ ,  $\gamma$ , bar_ $\sigma$ ,  $\mu_c$ ,  $\mu_d$  = sv_model
    I, J, K = len(hc_grid), len(hd_grid), len(z_grid)
    N = I * J * K

    A = compute_A(sv_model)
    # Make sure that a unique solution exists
    if test_stable:
        test_stability(A)

    # Compute v
    ones_array = np.ones(N)
    Id = np.identity(N)
    v = scipy.linalg.solve(Id - A, A @ ones_array)
    # Reshape into an array of the form v[i, j, k]
    v = np.reshape(v, (I, J, K))
    return v
```

Let's create an instance of the model and solve it.

```
sv_model = create_sv_model()
P, hc_grid, Q, hd_grid, R, z_grid,  $\beta$ ,  $\gamma$ , bar_ $\sigma$ ,  $\mu_c$ ,  $\mu_d$  = sv_model
```

Let's run it to compile.

```
start = time()
v = sv_pd_ratio(sv_model)
numpy_with_compile = time() - start
print("Numpy compile plus execution time = ", numpy_with_compile)
```

```
Numpy compile plus execution time = 0.7299137115478516
```

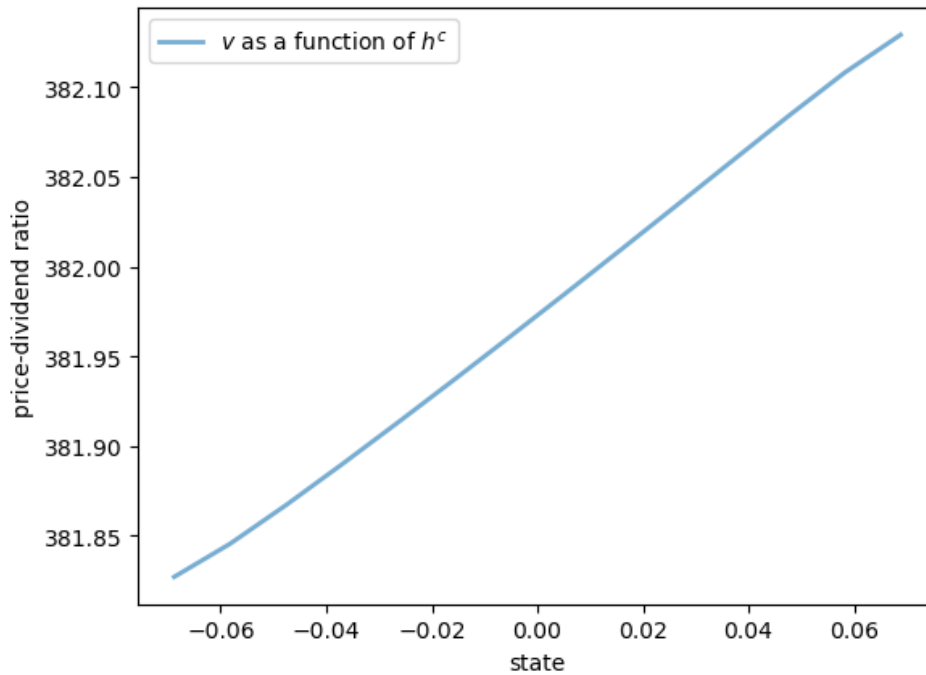
Let's run it again to remove the compile.

```
start = time()
v = sv_pd_ratio(sv_model)
numpy_without_compile = time() - start
print("Numpy execution time = ", numpy_without_compile)
```

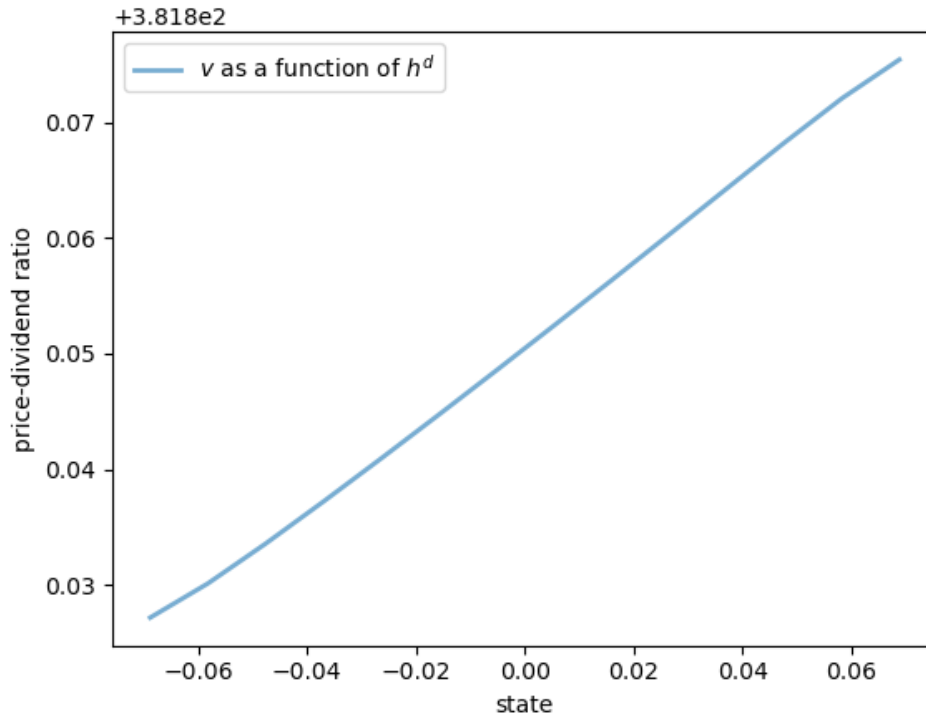
```
Numpy execution time = 0.2265474796295166
```

Here are some plots of the solution  $v$  along the three dimensions.

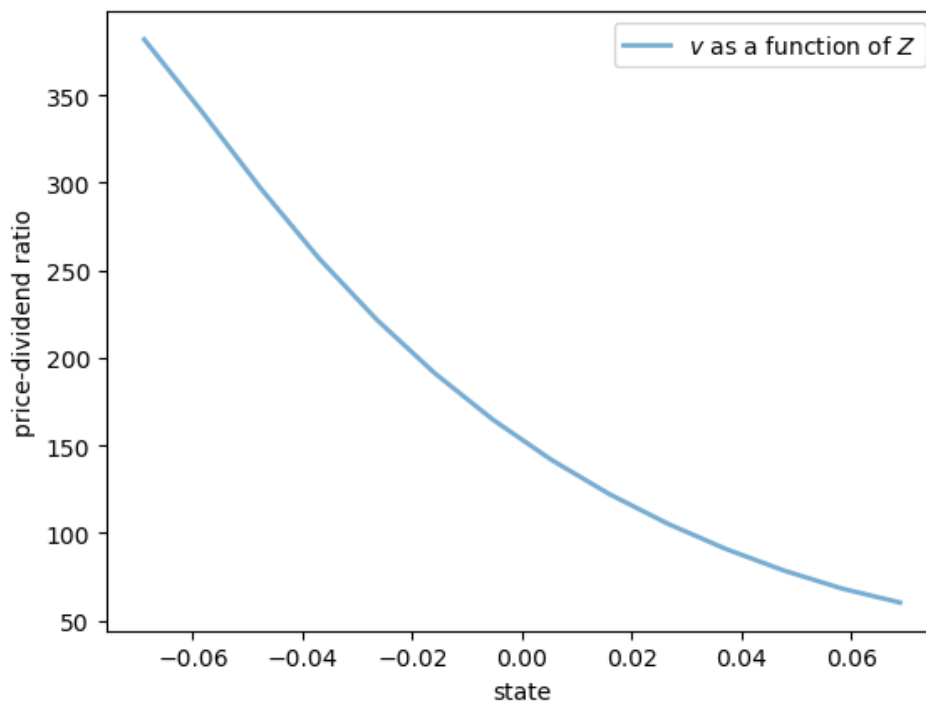
```
fig, ax = plt.subplots()
ax.plot(hc_grid, v[:, 0, 0], lw=2, alpha=0.6, label="$v$ as a function of $h^c$")
ax.set_ylabel("price-dividend ratio")
ax.set_xlabel("state")
ax.legend()
plt.show()
```



```
fig, ax = plt.subplots()
ax.plot(hd_grid, v[0, :, 0], lw=2, alpha=0.6, label="$v$ as a function of $h^d$")
ax.set_ylabel("price-dividend ratio")
ax.set_xlabel("state")
ax.legend()
plt.show()
```



```
fig, ax = plt.subplots()
ax.plot(z_grid, v[0, 0, :], lw=2, alpha=0.6, label="$v$ as a function of $Z$")
ax.set_ylabel("price-dividend ratio")
ax.set_xlabel("state")
ax.legend()
plt.show()
```





## 9.9 JAX Version

Now let's write a JAX version that is a simple transformation of the NumPy version.

(Below we will write a more efficient version using JAX's ability to work with linear operators.)

```
def create_sv_model_jax(sv_model):    # mean growth of dividends

    # Take the contents of a NumPy sv_model instance
    P, hc_grid, Q, hd_grid, R, z_grid,  $\beta$ ,  $\gamma$ , bar_ $\sigma$ ,  $\mu_c$ ,  $\mu_d$  = sv_model

    # Shift the arrays to the device (GPU if available)
    hc_grid, hd_grid, z_grid = map(jax.device_put, (hc_grid, hd_grid, z_grid))
    P, Q, R = map(jax.device_put, (P, Q, R))

    # Create a new instance and return it
    return SVMModel(P=P, hc_grid=hc_grid,
                    Q=Q, hd_grid=hd_grid,
                    R=R, z_grid=z_grid,
                     $\beta$ = $\beta$ ,  $\gamma$ = $\gamma$ , bar_ $\sigma$ =bar_ $\sigma$ ,  $\mu_c$ = $\mu_c$ ,  $\mu_d$ = $\mu_d$ )
```

Here's a function to compute  $A$ .

We include the extra argument shapes to help the compiler understand the size of the arrays.

This is important when we JIT-compile the function below.

```
def compute_A_jax(sv_model, shapes):
    # Set up
    P, hc_grid, Q, hd_grid, R, z_grid,  $\beta$ ,  $\gamma$ , bar_ $\sigma$ ,  $\mu_c$ ,  $\mu_d$  = sv_model
    I, J, K = shapes
    N = I * J * K
    # Reshape and broadcast over (i, j, k, i', j', k')
    hc = jnp.reshape(hc_grid, (I, 1, 1, 1, 1, 1))
    hd = jnp.reshape(hd_grid, (1, J, 1, 1, 1, 1))
    z = jnp.reshape(z_grid, (1, 1, K, 1, 1, 1))
    P = jnp.reshape(P, (I, 1, 1, I, 1, 1))
    Q = jnp.reshape(Q, (1, J, 1, 1, J, 1))
    R = jnp.reshape(R, (1, 1, K, 1, 1, K))
    # Compute A and then reshape to create a matrix
    a =  $\mu_d$  -  $\gamma$  *  $\mu_c$ 
    b = bar_ $\sigma$ **2 * (jnp.exp(2 * hd) +  $\gamma$ **2 * jnp.exp(2 * hc)) / 2
     $\kappa$  = jnp.exp(a + (1 -  $\gamma$ ) * z + b)
    A =  $\beta$  *  $\kappa$  * P * Q * R
    A = jnp.reshape(A, (N, N))
    return A
```

Here's the function that computes the solution.

```
def sv_pd_ratio_jax(sv_model, shapes):
    """
    Computes the price-dividend ratio of the asset for the stochastic volatility
    model.

    Parameters
    -----
    sv_model: an instance of Model
```

(continues on next page)

(continued from previous page)

```

    contains primitives

Returns
-----
v : array_like
    price-dividend ratio

"""
# unpack
P, hc_grid, Q, hd_grid, R, z_grid,  $\beta$ ,  $\gamma$ , bar_ $\sigma$ ,  $\mu_c$ ,  $\mu_d$  = sv_model
I, J, K = len(hc_grid), len(hd_grid), len(z_grid)
shapes = I, J, K
N = I * J * K

A = compute_A_jax(sv_model, shapes)

# Compute v, reshape and return
ones_array = jnp.ones(N)
Id = jnp.identity(N)
v = jax.scipy.linalg.solve(Id - A, A @ ones_array)
return jnp.reshape(v, (I, J, K))

```

Now let's target these functions for JIT-compilation, while using `static_argnums` to indicate that the function will need to be recompiled when shapes changes.

```

compute_A_jax = jax.jit(compute_A_jax, static_argnums=(1,))
sv_pd_ratio_jax = jax.jit(sv_pd_ratio_jax, static_argnums=(1,))

```

```

sv_model = create_sv_model()
sv_model_jax = create_sv_model_jax(sv_model)
P, hc_grid, Q, hd_grid, R, z_grid,  $\beta$ ,  $\gamma$ , bar_ $\sigma$ ,  $\mu_c$ ,  $\mu_d$  = sv_model_jax
shapes = len(hc_grid), len(hd_grid), len(z_grid)

```

Let's see how long it takes to run with compile time included.

```

start = time()
v_jax = sv_pd_ratio_jax(sv_model_jax, shapes).block_until_ready()
jnp_with_compile = time() - start
print("JAX compile plus execution time = ", jnp_with_compile)

```

```
JAX compile plus execution time = 0.4375894069671631
```

And now let's see without compile time.

```

start = time()
v_jax = sv_pd_ratio_jax(sv_model_jax, shapes).block_until_ready()
jnp_without_compile = time() - start
print("JAX execution time = ", jnp_without_compile)

```

```
JAX execution time = 0.09215307235717773
```

Here's the ratio of times:

```
jnp_without_compile / numpy_without_compile
```

```
0.40677156288774363
```

Let's check that the NumPy and JAX versions realize the same solution.

```
v = jax.device_put(v)
print(jnp.allclose(v, v_jax))
```

```
True
```

## 9.10 A memory-efficient JAX version

One problem with the code above is that we instantiate a matrix of size  $N = I \times J \times K$ .

This quickly becomes impossible as  $I, J, K$  increase.

Fortunately, JAX makes it possible to solve for the price-dividend ratio without instantiating this large matrix.

The first step is to think of  $A$  not as a matrix, but rather as the linear operator that transforms  $g$  into  $Ag$ .

```
def A(g, sv_model, shapes):
    # Set up
    P, hc_grid, Q, hd_grid, R, z_grid, beta, gamma, bar_sigma, mu_c, mu_d = sv_model
    I, J, K = shapes
    # Reshape and broadcast over (i, j, k, i', j', k')
    hc = jnp.reshape(hc_grid, (I, 1, 1, 1, 1, 1))
    hd = jnp.reshape(hd_grid, (1, J, 1, 1, 1, 1))
    z = jnp.reshape(z_grid, (1, 1, K, 1, 1, 1))
    P = jnp.reshape(P, (I, 1, 1, I, 1, 1))
    Q = jnp.reshape(Q, (1, J, 1, 1, J, 1))
    R = jnp.reshape(R, (1, 1, K, 1, 1, K))
    g = jnp.reshape(g, (1, 1, 1, I, J, K))
    a = mu_d - gamma * mu_c
    b = bar_sigma**2 * (jnp.exp(2 * hd) + gamma**2 * jnp.exp(2 * hc)) / 2
    kappa = jnp.exp(a + (1 - gamma) * z + b)
    A = beta * kappa * P * Q * R
    Ag = jnp.sum(A * g, axis=(3, 4, 5))
    return Ag
```

Now we write a version of the solution function for the price-dividend ratio that acts directly on the linear operator  $A$ .

```
def sv_pd_ratio_linop(sv_model, shapes):
    P, hc_grid, Q, hd_grid, R, z_grid, beta, gamma, bar_sigma, mu_c, mu_d = sv_model
    I, J, K = shapes

    ones_array = jnp.ones((I, J, K))
    # Set up the operator g -> (I - A) g
    J = lambda g: g - A(g, sv_model, shapes)
    # Solve v = (I - A)^{-1} A 1
    A1 = A(ones_array, sv_model, shapes)
    # Apply an iterative solver that works for linear operators
```

(continues on next page)

(continued from previous page)

```
v = jax.scipy.sparse.linalg.bicgstab(J, A1)[0]
return v
```

Let's target these functions for JIT compilation.

```
A = jax.jit(A, static_argnums=(2,))
sv_pd_ratio_linop = jax.jit(sv_pd_ratio_linop, static_argnums=(1,))
```

Let's time the solution with compile time included.

```
start = time()
v_jax_linop = sv_pd_ratio_linop(sv_model, shapes).block_until_ready()
jnp_linop_with_compile = time() - start
print("JAX compile plus execution time = ", jnp_linop_with_compile)
```

```
JAX compile plus execution time = 0.6296956539154053
```

And now let's see without compile time.

```
start = time()
v_jax_linop = sv_pd_ratio_linop(sv_model, shapes).block_until_ready()
jnp_linop_without_compile = time() - start
print("JAX execution time = ", jnp_linop_without_compile)
```

```
JAX execution time = 0.0056879520416259766
```

Let's verify the solution again:

```
print(jnp.allclose(v, v_jax_linop))
```

```
True
```

Here's the ratio of times between memory-efficient and direct version:

```
jnp_linop_without_compile / jnp_without_compile
```

```
0.061722869310096816
```

The speed is somewhat faster and, moreover, we can now work with much larger grids.

Here's a moderately large example, where the state space has 15,625 elements.

```
sv_model = create_sv_model(I=25, J=25, K=25)
sv_model_jax = create_sv_model_jax(sv_model)
P, hc_grid, Q, hd_grid, R, z_grid, beta, y, bar_sigma, mu_c, mu_d = sv_model_jax
shapes = len(hc_grid), len(hd_grid), len(z_grid)

%time _ = sv_pd_ratio_linop(sv_model_jax, shapes).block_until_ready()
%time _ = sv_pd_ratio_linop(sv_model_jax, shapes).block_until_ready()
```

```
CPU times: user 493 ms, sys: 2.59 ms, total: 496 ms  
Wall time: 692 ms  
CPU times: user 178 ms, sys: 0 ns, total: 178 ms  
Wall time: 178 ms
```

The solution is computed relatively quickly and without memory issues.

Readers will find that they can push these numbers further, although we refrain from doing so here.



## **Part IV**

# **Dynamic Programming**





## JOB SEARCH

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

In this lecture we study a basic infinite-horizon job search problem with Markov wage draws

---

**Note:** For background on infinite horizon job search see, e.g., [DP1](#).

---

The exercise at the end asks you to add risk-sensitive preferences and see how the main results change.

In addition to what’s in Anaconda, this lecture will need the following libraries:

```
!pip install quantecon
```

We use the following imports.

```
import matplotlib.pyplot as plt
import quantecon as qe
import jax
import jax.numpy as jnp
from collections import namedtuple

jax.config.update("jax_enable_x64", True)
```

## 10.1 Model

We study an elementary model where

- jobs are permanent
- unemployed workers receive current compensation  $c$
- the horizon is infinite
- an unemployment agent discounts the future via discount factor  $\beta \in (0, 1)$

### 10.1.1 Set up

At the start of each period, an unemployed worker receives wage offer  $W_t$ .

To build a wage offer process we consider the dynamics

$$W_{t+1} = \rho W_t + \nu Z_{t+1}$$

where  $(Z_t)_{t \geq 0}$  is IID and standard normal.

We then discretize this wage process using Tauchen's method to produce a stochastic matrix  $P$ .

Successive wage offers are drawn from  $P$ .

### 10.1.2 Rewards

Since jobs are permanent, the return to accepting wage offer  $w$  today is

$$w + \beta w + \beta^2 w + \dots = \frac{w}{1 - \beta}$$

The Bellman equation is

$$v(w) = \max \left\{ \frac{w}{1 - \beta}, c + \beta \sum_{w'} v(w') P(w, w') \right\}$$

We solve this model using value function iteration.

## 10.2 Code

Let's set up a namedtuple to store information needed to solve the model.

```
Model = namedtuple('Model', ('n', 'w_vals', 'P', 'β', 'c'))
```

The function below holds default values and populates the namedtuple.

```
def create_js_model(
    n=500,          # wage grid size
    ρ=0.9,          # wage persistence
    ν=0.2,          # wage volatility
    β=0.99,         # discount factor
    c=1.0,          # unemployment compensation
):
    "Creates an instance of the job search model with Markov wages."
    mc = qe.tauchen(n, ρ, ν)
    w_vals, P = jnp.exp(mc.state_values), jnp.array(mc.P)
    return Model(n, w_vals, P, β, c)
```

Let's test it:

```
model = create_js_model(β=0.98)
```

```
2024-08-12 03:59:23.172517: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages. (continues on next page)
```

(continued from previous page)

```
model.c
```

```
1.0
```

```
model.β
```

```
0.98
```

```
model.w_vals.mean()
```

```
Array(1.34861482, dtype=float64)
```

Here's the Bellman operator.

```
@jax.jit
def T(v, model):
    """
    The Bellman operator  $Tv = \max\{e, c + \beta E v\}$  with

     $e(w) = w / (1-\beta)$  and  $(Ev)(w) = E_w[v(W')]$ 

    """
    n, w_vals, P, β, c = model
    h = c + β * P @ v
    e = w_vals / (1 - β)

    return jnp.maximum(e, h)
```

The next function computes the optimal policy under the assumption that  $v$  is the value function.

The policy takes the form

$$\sigma(w) = \mathbf{1} \left\{ \frac{w}{1-\beta} \geq c + \beta \sum_{w'} v(w') P(w, w') \right\}$$

Here  $\mathbf{1}$  is an indicator function.

- $\sigma(w) = 1$  means stop
- $\sigma(w) = 0$  means continue.

```
@jax.jit
def get_greedy(v, model):
    "Get a v-greedy policy."
    n, w_vals, P, β, c = model
    e = w_vals / (1 - β)
    h = c + β * P @ v
    σ = jnp.where(e >= h, 1, 0)
    return σ
```

Here's a routine for value function iteration.

```
def vfi(model, max_iter=10_000, tol=1e-4):
    "Solve the infinite-horizon Markov job search model by VFI."
    print("Starting VFI iteration.")
    v = jnp.zeros_like(model.w_vals)    # Initial guess
    i = 0
    error = tol + 1

    while error > tol and i < max_iter:
        new_v = T(v, model)
        error = jnp.max(jnp.abs(new_v - v))
        i += 1
        v = new_v

    v_star = v
    sigma_star = get_greedy(v_star, model)
    return v_star, sigma_star
```

## 10.3 Computing the solution

Let's set up and solve the model.

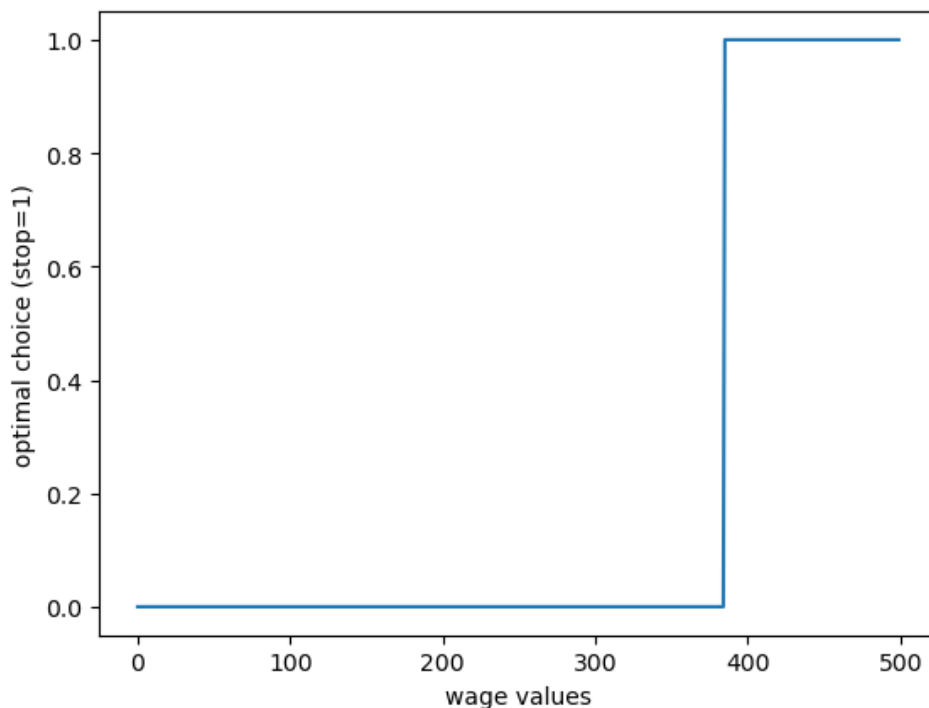
```
model = create_js_model()
n, w_vals, P, beta, c = model

v_star, sigma_star = vfi(model)
```

```
Starting VFI iteration.
```

Here's the optimal policy:

```
fig, ax = plt.subplots()
ax.plot(sigma_star)
ax.set_xlabel("wage values")
ax.set_ylabel("optimal choice (stop=1)")
plt.show()
```



We compute the reservation wage as the first  $w$  such that  $\sigma(w) = 1$ .

```
stop_indices = jnp.where( $\sigma\_star == 1$ )
stop_indices
```

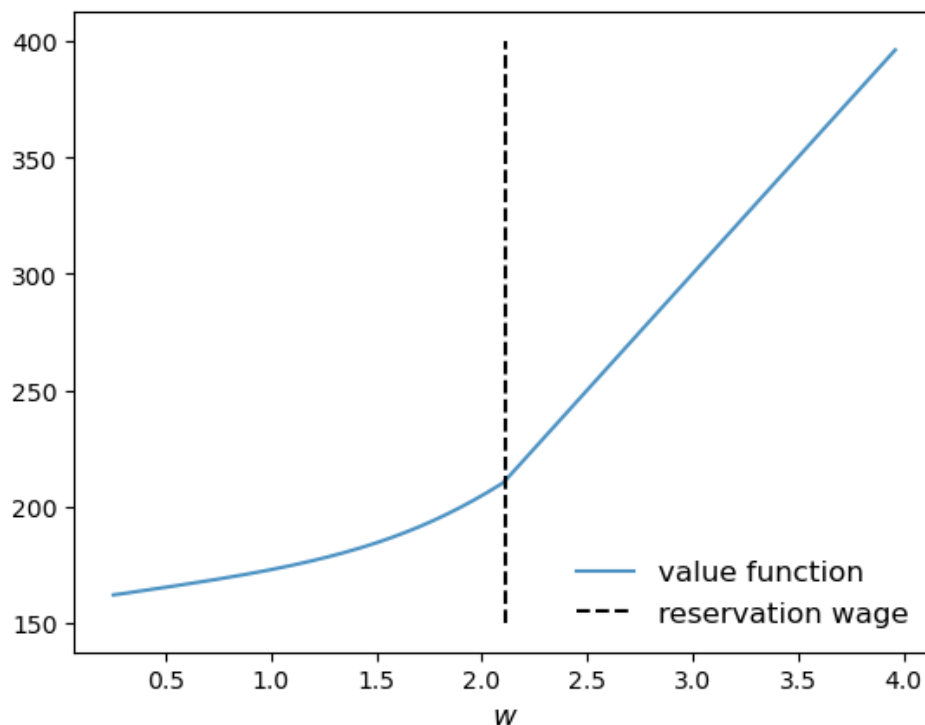
```
(Array([385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397,
        398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410,
        411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423,
        424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436,
        437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449,
        450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462,
        463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475,
        476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488,
        489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499],      dtype=int64),)
```

```
res_wage_index = min(stop_indices[0])
```

```
res_wage = w_vals[res_wage_index]
```

Here's a joint plot of the value function and the reservation wage.

```
fig, ax = plt.subplots()
ax.plot(w_vals, v_star, alpha=0.8, label="value function")
ax.vlines((res_wage,), 150, 400, 'k', ls='--', label="reservation wage")
ax.legend(frameon=False, fontsize=12, loc="lower right")
ax.set_xlabel("$w$", fontsize=12)
plt.show()
```



## 10.4 Exercise

### Exercise 10.4.1

In the setting above, the agent is risk-neutral vis-a-vis future utility risk.

Now solve the same problem but this time assuming that the agent has risk-sensitive preferences, which are a type of nonlinear recursive preferences.

The Bellman equation becomes

$$v(w) = \max \left\{ \frac{w}{1-\beta}, c + \frac{\beta}{\theta} \ln \left[ \sum_{w'} \exp(\theta v(w')) P(w, w') \right] \right\}$$

When  $\theta < 0$  the agent is risk averse.

Solve the model when  $\theta = -0.1$  and compare your result to the risk neutral case.

Try to interpret your result.

You can start with the following code:

```
RiskModel = namedtuple('Model', ('n', 'w_vals', 'P', 'β', 'c', 'θ'))

def create_risk_sensitive_js_model(
    n=500,          # wage grid size
    ρ=0.9,          # wage persistence
    v=0.2,          # wage volatility
    β=0.99,         # discount factor
```

(continues on next page)

(continued from previous page)

```

        c=1.0,          # unemployment compensation
        θ=-0.1         # risk parameter
    ):
        "Creates an instance of the job search model with Markov wages."
        mc = qe.tauchen(n, ρ, v)
        w_vals, P = jnp.exp(mc.state_values), mc.P
        P = jnp.array(P)
        return RiskModel(n, w_vals, P, β, c, θ)

```

Now you need to modify `T` and `get_greedy` and then run value function iteration again.

### Solution to Exercise 10.4.1

```
RiskModel = namedtuple('Model', ('n', 'w_vals', 'P', 'β', 'c', 'θ'))
```

```

def create_risk_sensitive_js_model(
    n=500,          # wage grid size
    ρ=0.9,          # wage persistence
    v=0.2,          # wage volatility
    β=0.99,         # discount factor
    c=1.0,          # unemployment compensation
    θ=-0.1         # risk parameter
):
    "Creates an instance of the job search model with Markov wages."
    mc = qe.tauchen(n, ρ, v)
    w_vals, P = jnp.exp(mc.state_values), mc.P
    P = jnp.array(P)
    return RiskModel(n, w_vals, P, β, c, θ)

```

```

@jax.jit
def T_rs(v, model):
    """
    The Bellman operator  $Tv = \max\{e, c + \beta R v\}$  with

    
$$e(w) = w / (1 - \beta) \text{ and}$$

    
$$(Rv)(w) = (1/\theta) \ln\{E_w[\exp(\theta v(W'))]\}$$


    """
    n, w_vals, P, β, c, θ = model
    h = c + (β / θ) * jnp.log(P @ (jnp.exp(θ * v)))
    e = w_vals / (1 - β)

    return jnp.maximum(e, h)

```

```

@jax.jit
def get_greedy_rs(v, model):
    "Get a v-greedy policy."
    n, w_vals, P, β, c, θ = model
    e = w_vals / (1 - β)
    h = c + (β / θ) * jnp.log(P @ (jnp.exp(θ * v)))
    σ = jnp.where(e >= h, 1, 0)

```

(continues on next page)

(continued from previous page)

```

    return  $\sigma$ 

def vfi(model, max_iter=10_000, tol=1e-4):
    "Solve the infinite-horizon Markov job search model by VFI."
    print("Starting VFI iteration.")
    v = jnp.zeros_like(model.w_vals)    # Initial guess
    i = 0
    error = tol + 1

    while error > tol and i < max_iter:
        new_v = T_rs(v, model)
        error = jnp.max(jnp.abs(new_v - v))
        i += 1
        v = new_v

    v_star = v
     $\sigma$ _star = get_greedy_rs(v_star, model)
    return v_star,  $\sigma$ _star

model_rs = create_risk_sensitive_js_model()

n, w_vals, P,  $\beta$ , c,  $\theta$  = model_rs

v_star_rs,  $\sigma$ _star_rs = vfi(model_rs)

```

```
Starting VFI iteration.
```

Let's plot the results together with the original risk neutral case and see what we get.

```

stop_indices = jnp.where( $\sigma$ _star_rs == 1)
res_wage_index = min(stop_indices[0])
res_wage_rs = w_vals[res_wage_index]

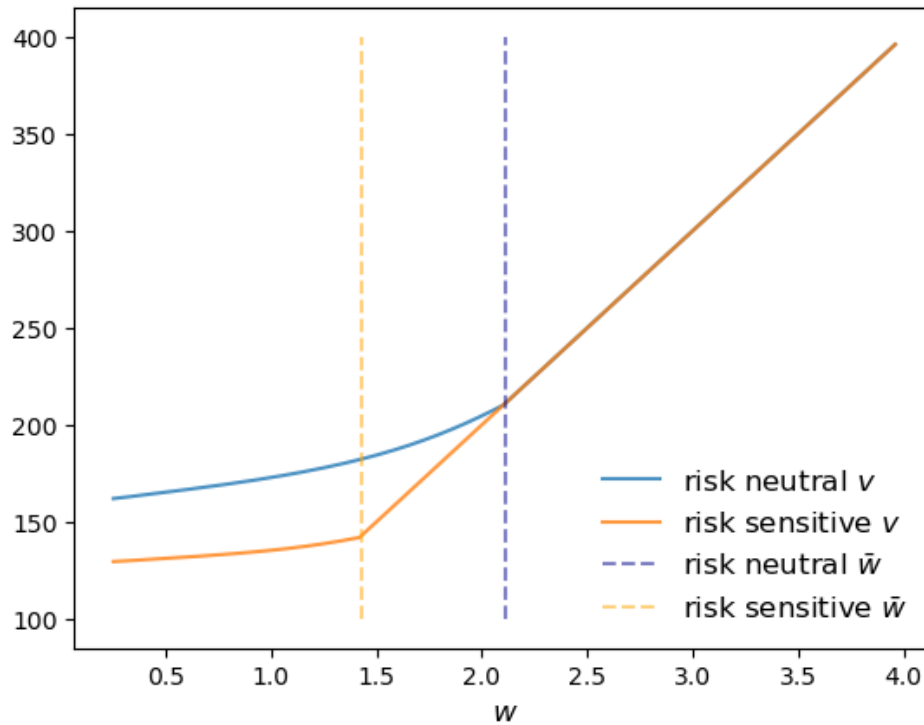
```

```

fig, ax = plt.subplots()
ax.plot(w_vals, v_star, alpha=0.8, label="risk neutral $v$")
ax.plot(w_vals, v_star_rs, alpha=0.8, label="risk sensitive $v$")
ax.vlines((res_wage,), 100, 400, ls='--', color='darkblue',
          alpha=0.5, label=r"risk neutral $\bar{w}$")
ax.vlines((res_wage_rs), 100, 400, ls='--', color='orange',
          alpha=0.5, label=r"risk sensitive $\bar{w}$")
ax.legend(frameon=False, fontsize=12, loc="lower right")
ax.set_xlabel("$w$", fontsize=12)
plt.show()

```





The figure shows that the reservation wage under risk sensitive preferences (RS  $\bar{w}$ ) shifts down.

This makes sense – the agent does not like risk and hence is more inclined to accept the current offer, even when it's lower.



## OPTIMAL SAVINGS I: VALUE FUNCTION ITERATION

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

In addition to JAX and Anaconda, this lecture will need the following libraries:

```
!pip install quantecon
```

We will use the following imports:

```
import quantecon as qe
import numpy as np
import jax
import jax.numpy as jnp
import matplotlib.pyplot as plt
from collections import namedtuple
from time import time
```

Let’s check the GPU we are running

```
!nvidia-smi
```

```
Mon Aug 12 04:02:56 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.0
| 5                                     |
+-----+-----+-----+-----+
| GPU   Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC   |                                           |                 |        |
| Fan   Temp   Perf              Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |                                           |              |          |
| MIG M. |                                           |              |          |
+=====+-----+-----+-----+=====+
```

(continues on next page)

(continued from previous page)

```

| 0 Tesla T4 On | 00000001:00:00.0 Off |
↪ Off |
| N/A 47C P8 10W / 70W | 2MiB / 16384MiB | 0%
↪ Default |
|
| N/A |
+-----+
↪-----+
+-----+
| Processes:
↪ |
| GPU GI CI PID Type Process name GPU
↪Memory |
| ID ID Usage
↪ |
|=====|
| No running processes found
↪ |
+-----+
↪-----+

```

We'll use 64 bit floats to gain extra precision.

```
jax.config.update("jax_enable_x64", True)
```

## 11.1 Overview

We consider an optimal savings problem with CRRA utility and budget constraint

$$W_{t+1} + C_t \leq RW_t + Y_t$$

where

- $C_t$  is consumption and  $C_t \geq 0$ ,
- $W_t$  is wealth and  $W_t \geq 0$ ,
- $R > 0$  is a gross rate of return, and
- $(Y_t)$  is labor income.

We assume below that labor income is a discretized AR(1) process.

The Bellman equation is

$$v(w) = \max_{0 \leq w' \leq Rw+y} \left\{ u(Rw + y - w') + \beta \sum_{y'} v(w', y') Q(y, y') \right\}$$

where

$$u(c) = \frac{c^{1-\gamma}}{1-\gamma}$$

In the code we use the function

$$B((w, y), w', v) = u(Rw + y - w') + \beta \sum_{y'} v(w', y') Q(y, y').$$

the encapsulate the right hand side of the Bellman equation.

## 11.2 Starting with NumPy

Let's start with a standard NumPy version running on the CPU.

Starting with this traditional approach will allow us to record the speed gain associated with switching to JAX.

(NumPy operations are similar to MATLAB operations, so this also serves as a rough comparison with MATLAB.)

### 11.2.1 Functions and operators

The following function contains default parameters and returns tuples that contain the key computational components of the model.

```
def create_consumption_model(R=1.01,                # Gross interest rate
                             β=0.98,                # Discount factor
                             γ=2,                   # CRRA parameter
                             w_min=0.01,            # Min wealth
                             w_max=5.0,             # Max wealth
                             w_size=150,            # Grid side
                             ρ=0.9, v=0.1, y_size=100): # Income parameters
    """
    A function that takes in parameters and returns parameters and grids
    for the optimal savings problem.
    """
    # Build grids and transition probabilities
    w_grid = np.linspace(w_min, w_max, w_size)
    mc = qe.tauchen(n=y_size, rho=ρ, sigma=v)
    y_grid, Q = np.exp(mc.state_values), mc.P
    # Pack and return
    params = β, R, γ
    sizes = w_size, y_size
    arrays = w_grid, y_grid, Q
    return params, sizes, arrays
```

(The function returns sizes of arrays because we use them later to help compile functions in JAX.)

To produce efficient NumPy code, we will use a vectorized approach.

The first step is to create the right hand side of the Bellman equation as a multi-dimensional array with dimensions over all states and controls.

```
def B(v, params, sizes, arrays):
    """
    A vectorized version of the right-hand side of the Bellman equation
    (before maximization), which is a 3D array representing

    B(w, y, w') = u(Rw + y - w') + β Σ_{y'} v(w', y') Q(y, y')
```

(continues on next page)

(continued from previous page)

```

for all (w, y, w').
"""

# Unpack
β, R, γ = params
w_size, y_size = sizes
w_grid, y_grid, Q = arrays

# Compute current rewards r(w, y, wp) as array r[i, j, ip]
w = np.reshape(w_grid, (w_size, 1, 1))    # w[i]    -> w[i, j, ip]
y = np.reshape(y_grid, (1, y_size, 1))    # z[j]    -> z[i, j, ip]
wp = np.reshape(w_grid, (1, 1, w_size))   # wp[ip] -> wp[i, j, ip]
c = R * w + y - wp

# Calculate continuation rewards at all combinations of (w, y, wp)
v = np.reshape(v, (1, 1, w_size, y_size)) # v[ip, jp] -> v[i, j, ip, jp]
Q = np.reshape(Q, (1, y_size, 1, y_size)) # Q[j, jp] -> Q[i, j, ip, jp]
EV = np.sum(v * Q, axis=3)                # sum over last index jp

# Compute the right-hand side of the Bellman equation
return np.where(c > 0, c**(1-γ)/(1-γ) + β * EV, -np.inf)

```

Here are two functions we need for value function iteration.

The first is the Bellman operator.

The second computes a  $v$ -greedy policy given  $v$  (i.e., the policy that maximizes the right-hand side of the Bellman equation.)

```

def T(v, params, sizes, arrays):
    "The Bellman operator."
    return np.max(B(v, params, sizes, arrays), axis=2)

def get_greedy(v, params, sizes, arrays):
    "Computes a v-greedy policy, returned as a set of indices."
    return np.argmax(B(v, params, sizes, arrays), axis=2)

```

## 11.2.2 Value function iteration

Here's a routine that performs value function iteration.

```

def value_function_iteration(model, max_iter=10_000, tol=1e-5):
    params, sizes, arrays = model
    v = np.zeros(sizes)
    i, error = 0, tol + 1
    while error > tol and i < max_iter:
        v_new = T(v, params, sizes, arrays)
        error = np.max(np.abs(v_new - v))
        i += 1
        v = v_new
    return v, get_greedy(v, params, sizes, arrays)

```

Now we create an instance, unpack it, and test how long it takes to solve the model.

```

model = create_consumption_model()
# Unpack
params, sizes, arrays = model
 $\beta$ , R,  $\gamma$  = params
w_size, y_size = sizes
w_grid, y_grid, Q = arrays

```

```

print("Starting VFI.")
start = time()
v_star,  $\sigma$ _star = value_function_iteration(model)
numpy_with_compile = time() - start
print(f"VFI completed in {numpy_with_compile} seconds.")

```

Starting VFI.

VFI completed in 11.466776371002197 seconds.

Let's run it again to eliminate compile time.

```

start = time()
v_star,  $\sigma$ _star = value_function_iteration(model)
numpy_without_compile = time() - start
print(f"VFI completed in {numpy_without_compile} seconds.")

```

VFI completed in 11.510318756103516 seconds.

Here's a plot of the policy function.

```

fig, ax = plt.subplots()
ax.plot(w_grid, w_grid, "k--", label="45")
ax.plot(w_grid, w_grid[ $\sigma$ _star[:, 1]], label="$\\sigma^{*}(\cdot, y_1)$")
ax.plot(w_grid, w_grid[ $\sigma$ _star[:, -1]], label="$\\sigma^{*}(\cdot, y_N)$")
ax.legend()
plt.show()

```

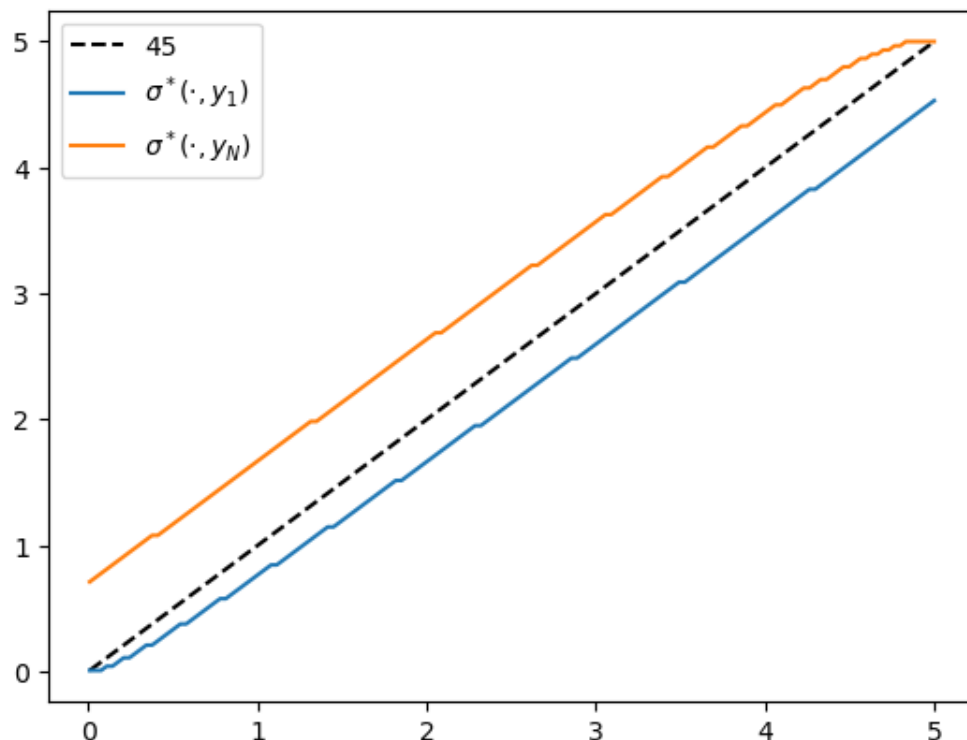


Fig. 11.1: Policy function

## 11.3 Switching to JAX

To switch over to JAX, we change `np` to `jnp` throughout and add some `jax.jit` requests.

### 11.3.1 Functions and operators

We redefine `create_consumption_model` to produce JAX arrays.

```
def create_consumption_model(R=1.01,                # Gross interest rate
                             beta=0.98,            # Discount factor
                             gamma=2,              # CRRA parameter
                             w_min=0.01,           # Min wealth
                             w_max=5.0,            # Max wealth
                             w_size=150,           # Grid size
                             rho=0.9, v=0.1, y_size=100): # Income parameters
    """
    A function that takes in parameters and returns parameters and grids
    for the optimal savings problem.
    """
    w_grid = jnp.linspace(w_min, w_max, w_size)
    mc = qe.tauchen(n=y_size, rho=rho, sigma=v)
    y_grid, Q = jnp.exp(mc.state_values), jax.device_put(mc.P)
    sizes = w_size, y_size
    return (beta, R, gamma), sizes, (w_grid, y_grid, Q)
```

The right hand side of the Bellman equation is the same as the NumPy version after switching `np` to `jnp`.

```
def B(v, params, sizes, arrays):
```

**146 Chapter 11. Optimal Savings I: Value Function Iteration**

"""  
A vectorized version of the right-hand side of the Bellman equation  
(before maximization), which is a 3D array representing

$B(w, y, w') = u(Rw + y - w') + \beta \sum_{y'} v(w', y') Q(y, y')$



Some readers might be concerned that we are creating high dimensional arrays, leading to inefficiency.

Could they be avoided by more careful vectorization?

In fact this is not necessary: this function will be JIT-compiled by JAX, and the JIT compiler will optimize compiled code to minimize memory use.

```
B = jax.jit(B, static_argnums=(2,))
```

In the call above, we indicate to the compiler that `sizes` is static, so the compiler can parallelize optimally while taking array sizes as fixed.

The Bellman operator  $T$  can be implemented by

```
def T(v, params, sizes, arrays):
    "The Bellman operator."
    return jnp.max(B(v, params, sizes, arrays), axis=2)

T = jax.jit(T, static_argnums=(2,))
```

The next function computes a  $v$ -greedy policy given  $v$  (i.e., the policy that maximizes the right-hand side of the Bellman equation.)

```
def get_greedy(v, params, sizes, arrays):
    "Computes a v-greedy policy, returned as a set of indices."
    return jnp.argmax(B(v, params, sizes, arrays), axis=2)

get_greedy = jax.jit(get_greedy, static_argnums=(2,))
```

## 11.3.2 Successive approximation

Now we define a solver that implements VFI.

We could use the one we built for NumPy above, after changing `np` to `jnp`.

Alternatively, we can push a bit harder and write a compiled version using `jax.lax.while_loop`.

This will give us just a bit more speed.

The first step is to write a compiled successive approximation routine that performs fixed point iteration on some given function  $T$ .

```
def successive_approx_jax(T,                # Operator (callable)
                          x_0,              # Initial condition
                          tolerance=1e-6,  # Error tolerance
                          max_iter=10_000): # Max iteration bound
    def body_fun(k_x_err):
        k, x, error = k_x_err
        x_new = T(x)
        error = jnp.max(jnp.abs(x_new - x))
        return k + 1, x_new, error

    def cond_fun(k_x_err):
        k, x, error = k_x_err
        return jnp.logical_and(error > tolerance, k < max_iter)

    k, x, error = jax.lax.while_loop(cond_fun, body_fun,
```

(continues on next page)

(continued from previous page)

```

                                (1, x_0, tolerance + 1))
    return x

successive_approx_jax = \
    jax.jit(successive_approx_jax, static_argnums=(0,))

```

Our value function iteration routine calls `successive_approx_jax` while passing in the Bellman operator.

```

def value_function_iteration(model, tol=1e-5):
    params, sizes, arrays = model
    vz = jnp.zeros(sizes)
    _T = lambda v: T(v, params, sizes, arrays)
    v_star = successive_approx_jax(_T, vz, tolerance=tol)
    return v_star, get_greedy(v_star, params, sizes, arrays)

```

### 11.3.3 Timing

Let's create an instance and unpack it.

```

model = create_consumption_model()
# Unpack
params, sizes, arrays = model
β, R, γ = params
w_size, y_size = sizes
w_grid, y_grid, Q = arrays

```

```

2024-08-12 04:03:20.795231: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.

```

Let's see how long it takes to solve this model.

```

print("Starting VFI using vectorization.")
start = time()
v_star_jax, σ_star_jax = value_function_iteration(model)
jax_with_compile = time() - start
print(f"VFI completed in {jax_with_compile} seconds.")

```

```
Starting VFI using vectorization.
```

```
VFI completed in 0.7179312705993652 seconds.
```

Let's run it again to eliminate compile time.

```

start = time()
v_star_jax, σ_star_jax = value_function_iteration(model)
jax_without_compile = time() - start
print(f"VFI completed in {jax_without_compile} seconds.")

```

```
VFI completed in 0.31941699981689453 seconds.
```

The relative speed gain is

```
print(f"Relative speed gain = {numpy_without_compile / jax_without_compile}")
```

```
Relative speed gain = 36.035398124401
```

This is an impressive speed up and in fact we can do better still by switching to alternative algorithms that are better suited to parallelization.

These algorithms are discussed in a [separate lecture](#).

## 11.4 Switching to vmap

Before we discuss alternative algorithms, let's take another look at value function iteration.

For this simple optimal savings problem, direct vectorization is relatively easy.

In particular, it's straightforward to express the right hand side of the Bellman equation as an array that stores evaluations of the function at every state and control.

For more complex models direct vectorization can be much harder.

For this reason, it helps to have another approach to fast JAX implementations up our sleeves.

Here's a version that

1. writes the right hand side of the Bellman operator as a function of individual states and controls, and
2. applies `jax.vmap` on the outside to achieve a parallelized solution.

First let's rewrite B

```
def B(v, params, arrays, i, j, ip):
    """
    The right-hand side of the Bellman equation before maximization, which takes
    the form

        
$$B(w, y, w') = u(Rw + y - w') + \beta \sum_{y'} v(w', y') Q(y, y')$$


    The indices are (i, j, ip) -> (w, y, w').
    """
    beta, R, gamma = params
    w_grid, y_grid, Q = arrays
    w, y, wp = w_grid[i], y_grid[j], w_grid[ip]
    c = R * w + y - wp
    EV = jnp.sum(v[ip, :] * Q[j, :])
    return jnp.where(c > 0, c**(1-gamma)/(1-gamma) + beta * EV, -jnp.inf)
```

Now we successively apply vmap to simulate nested loops.

```
B_1 = jax.vmap(B, in_axes=(None, None, None, None, None, 0))
B_2 = jax.vmap(B_1, in_axes=(None, None, None, None, 0, None))
B_vmap = jax.vmap(B_2, in_axes=(None, None, None, 0, None, None))
```

Here's the Bellman operator and the `get_greedy` functions for the vmap case.

```
def T_vmap(v, params, sizes, arrays):
    "The Bellman operator."
    w_size, y_size = sizes
    w_indices, y_indices = jnp.arange(w_size), jnp.arange(y_size)
    B_values = B_vmap(v, params, arrays, w_indices, y_indices, w_indices)
    return jnp.max(B_values, axis=-1)

T_vmap = jax.jit(T_vmap, static_argnums=(2,))

def get_greedy_vmap(v, params, sizes, arrays):
    "Computes a v-greedy policy, returned as a set of indices."
    w_size, y_size = sizes
    w_indices, y_indices = jnp.arange(w_size), jnp.arange(y_size)
    B_values = B_vmap(v, params, arrays, w_indices, y_indices, w_indices)
    return jnp.argmax(B_values, axis=-1)

get_greedy_vmap = jax.jit(get_greedy_vmap, static_argnums=(2,))
```

Here's the iteration routine.

```
def value_iteration_vmap(model, tol=1e-5):
    params, sizes, arrays = model
    vz = jnp.zeros(sizes)
    _T = lambda v: T_vmap(v, params, sizes, arrays)
    v_star = successive_approx_jax(_T, vz, tolerance=tol)
    return v_star, get_greedy(v_star, params, sizes, arrays)
```

Let's see how long it takes to solve the model using the vmap method.

```
print("Starting VFI using vmap.")
start = time()
v_star_vmap, sigma_star_vmap = value_iteration_vmap(model)
jax_vmap_with_compile = time() - start
print(f"VFI completed in {jax_vmap_with_compile} seconds.")
```

Starting VFI using vmap.

VFI completed in 0.4236183166503906 seconds.

Let's run it again to get rid of compile time.

```
start = time()
v_star_vmap, sigma_star_vmap = value_iteration_vmap(model)
jax_vmap_without_compile = time() - start
print(f"VFI completed in {jax_vmap_without_compile} seconds.")
```

VFI completed in 0.3306581974029541 seconds.

We need to make sure that we got the same result.

```
print(jnp.allclose(v_star_vmap, v_star_jax))
print(jnp.allclose(sigma_star_vmap, sigma_star_jax))
```

```
True
True
```

Here's the speed gain associated with switching from the NumPy version to JAX with `vmap`:

```
print(f"Relative speed = {numpy_without_compile/jax_vmap_without_compile}")
```

```
Relative speed = 34.810323308200196
```

And here's the comparison with the first JAX implementation (which used direct vectorization).

```
print(f"Relative speed = {jax_without_compile / jax_vmap_without_compile}")
```

```
Relative speed = 0.9660035720440326
```

The execution times for the two JAX versions are relatively similar.

However, as emphasized above, having a second method up our sleeves (i.e, the `vmap` approach) will be helpful when confronting dynamic programs with more sophisticated Bellman equations.



## OPTIMAL SAVINGS II: ALTERNATIVE ALGORITHMS

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

In *Optimal Savings I: Value Function Iteration* we solved a simple version of the household optimal savings problem via value function iteration (VFI) using JAX.

In this lecture we tackle exactly the same problem while adding in two alternative algorithms:

- optimistic policy iteration (OPI) and
- Howard policy iteration (HPI).

We will see that both of these algorithms outperform traditional VFI.

One reason for this is that the algorithms have good convergence properties.

Another is that one of them, HPI, is particularly well suited to pairing with JAX.

The reason is that HPI uses a relatively small number of computationally expensive steps, whereas VFI uses a longer sequence of small steps.

In other words, VFI is inherently more sequential than HPI, and sequential routines are hard to parallelize.

By comparison, HPI is less sequential – the small number of computationally intensive steps can be effectively parallelized by JAX.

This is particularly valuable when the underlying hardware includes a GPU.

Details on VFI, HPI and OPI can be found in [this book](#), for which a PDF is freely available.

Here we assume readers have some knowledge of the algorithms and focus on computation.

For the details of the savings model, readers can refer to *Optimal Savings I: Value Function Iteration*.

In addition to what’s in Anaconda, this lecture will need the following libraries:

```
!pip install quantecon
```

We will use the following imports:

```
import quantecon as qe
import jax
import jax.numpy as jnp
import matplotlib.pyplot as plt
from time import time
```

Let's check the GPU we are running.

```
!nvidia-smi
```

```
Mon Aug 12 04:03:28 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03   CUDA Version: 12.
| 5 |
+-----+-----+-----+
| GPU Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+=====+=====+=====+
|    0   Tesla T4                          On   | 00000001:00:00.0 Off  |
| Off |
| N/A   42C    P8              10W / 70W |      2MiB / 16384MiB |      0%   |
| Default |
|
| N/A |
+-----+-----+-----+

+-----+
| Processes:
|
| GPU   GI    CI        PID   Type   Process name                      GPU
| Memory |
|      ID    ID
|
|
+=====+=====+=====+
| No running processes found
|
+-----+
```

We'll use 64 bit floats to gain extra precision.

```
jax.config.update("jax_enable_x64", True)
```



## 12.1 Model primitives

First we define a model that stores parameters and grids.

The following code is repeated from *Optimal Savings I: Value Function Iteration*.

```
def create_consumption_model(R=1.01,                # Gross interest rate
                             β=0.98,               # Discount factor
                             γ=2,                  # CRRA parameter
                             w_min=0.01,           # Min wealth
                             w_max=5.0,           # Max wealth
                             w_size=150,           # Grid side
                             ρ=0.9, v=0.1, y_size=100): # Income parameters
    """
    A function that takes in parameters and returns parameters and grids
    for the optimal savings problem.
    """
    # Build grids and transition probabilities
    w_grid = jnp.linspace(w_min, w_max, w_size)
    mc = qe.tauchen(n=y_size, rho=ρ, sigma=v)
    y_grid, Q = jnp.exp(mc.state_values), mc.P
    # Pack and return
    params = β, R, γ
    sizes = w_size, y_size
    arrays = w_grid, y_grid, jnp.array(Q)
    return params, sizes, arrays
```

Here's the right hand side of the Bellman equation:

```
def _B(v, params, arrays, i, j, ip):
    """
    The right-hand side of the Bellman equation before maximization, which takes
    the form

        
$$B(w, y, w') = u(Rw + y - w') + \beta \sum_{y'} v(w', y') Q(y, y')$$


    The indices are (i, j, ip) -> (w, y, w').
    """
    β, R, γ = params
    w_grid, y_grid, Q = arrays
    w, y, wp = w_grid[i], y_grid[j], w_grid[ip]
    c = R * w + y - wp
    EV = jnp.sum(v[ip, :] * Q[j, :])
    return jnp.where(c > 0, c**(1-γ)/(1-γ) + β * EV, -jnp.inf)
```

Now we successively apply vmap to vectorize  $B$  by simulating nested loops.

```
B_1 = jax.vmap(_B, in_axes=(None, None, None, None, None, 0))
B_2 = jax.vmap(B_1, in_axes=(None, None, None, None, 0, None))
B_vmap = jax.vmap(B_2, in_axes=(None, None, None, 0, None, None))
```

Here's a fully vectorized version of  $B$ .

```
def B(v, params, sizes, arrays):
    w_size, y_size = sizes
    w_indices, y_indices = jnp.arange(w_size), jnp.arange(y_size)
```

(continues on next page)

(continued from previous page)

```

    return B_vmap(v, params, arrays, w_indices, y_indices, w_indices)

B = jax.jit(B, static_argnums=(2,))

```

## 12.2 Operators

Here's the Bellman operator  $T$

```

def T(v, params, sizes, arrays):
    "The Bellman operator."
    return jnp.max(B(v, params, sizes, arrays), axis=-1)

T = jax.jit(T, static_argnums=(2,))

```

The next function computes a  $v$ -greedy policy given  $v$

```

def get_greedy(v, params, sizes, arrays):
    "Computes a v-greedy policy, returned as a set of indices."
    return jnp.argmax(B(v, params, sizes, arrays), axis=-1)

get_greedy = jax.jit(get_greedy, static_argnums=(2,))

```

We define a function to compute the current rewards  $r_\sigma$  given policy  $\sigma$ , which is defined as the vector

$$r_\sigma(w, y) := r(w, y, \sigma(w, y))$$

```

def _compute_r_sigma(sigma, params, arrays, i, j):
    """
    With indices (i, j) -> (w, y) and wp = sigma[i, j], compute

        r_sigma[i, j] = u(Rw + y - wp)

    which gives current rewards under policy sigma.
    """

    # Unpack model
    beta, R, gamma = params
    w_grid, y_grid, Q = arrays
    # Compute r_sigma[i, j]
    w, y, wp = w_grid[i], y_grid[j], w_grid[sigma[i, j]]
    c = R * w + y - wp
    r_sigma = c ** (1 - gamma) / (1 - gamma)

    return r_sigma

```

Now we successively apply vmap to simulate nested loops.

```

r_1 = jax.vmap(_compute_r_sigma, in_axes=(None, None, None, None, 0))
r_sigma_vmap = jax.vmap(r_1, in_axes=(None, None, None, 0, None))

```

Here's a fully vectorized version of  $r_\sigma$ .

```
def compute_r_σ(σ, params, sizes, arrays):
    w_size, y_size = sizes
    w_indices, y_indices = jnp.arange(w_size), jnp.arange(y_size)
    return r_σ_vmap(σ, params, arrays, w_indices, y_indices)

compute_r_σ = jax.jit(compute_r_σ, static_argnums=(2,))
```

Now we define the policy operator  $T_\sigma$  going through similar steps

```
def _T_σ(v, σ, params, arrays, i, j):
    "The σ-policy operator."

    # Unpack model
    β, R, γ = params
    w_grid, y_grid, Q = arrays

    r_σ = _compute_r_σ(σ, params, arrays, i, j)
    # Calculate the expected sum  $\sum_{jp} v[\sigma[i, j], jp] * Q[i, j, jp]$ 
    EV = jnp.sum(v[σ[i, j], :] * Q[j, :])

    return r_σ + β * EV

T_1 = jax.vmap(_T_σ, in_axes=(None, None, None, None, None, 0))
T_σ_vmap = jax.vmap(T_1, in_axes=(None, None, None, None, 0, None))

def T_σ(v, σ, params, sizes, arrays):
    w_size, y_size = sizes
    w_indices, y_indices = jnp.arange(w_size), jnp.arange(y_size)
    return T_σ_vmap(v, σ, params, arrays, w_indices, y_indices)

T_σ = jax.jit(T_σ, static_argnums=(3,))
```

The function below computes the value  $v_\sigma$  of following policy  $\sigma$ .

This lifetime value is a function  $v_\sigma$  that satisfies

$$v_\sigma(w, y) = r_\sigma(w, y) + \beta \sum_{y'} v_\sigma(\sigma(w, y), y') Q(y, y')$$

We wish to solve this equation for  $v_\sigma$ .

Suppose we define the linear operator  $L_\sigma$  by

$$(L_\sigma v)(w, y) = v(w, y) - \beta \sum_{y'} v(\sigma(w, y), y') Q(y, y')$$

With this notation, the problem is to solve for  $v$  via

$$(L_\sigma v)(w, y) = r_\sigma(w, y)$$

In vector for this is  $L_\sigma v = r_\sigma$ , which tells us that the function we seek is

$$v_\sigma = L_\sigma^{-1} r_\sigma$$

JAX allows us to solve linear systems defined in terms of operators; the first step is to define the function  $L_\sigma$ .

```

def _L_σ(v, σ, params, arrays, i, j):
    """
    Here we set up the linear map  $v \rightarrow L_\sigma v$ , where

         $(L_\sigma v)(w, y) = v(w, y) - \beta \sum_{y'} v(\sigma(w, y), y') Q(y, y')$ 

    """
    # Unpack
    β, R, γ = params
    w_grid, y_grid, Q = arrays
    # Compute and return  $v[i, j] - \beta \sum_{jp} v[\sigma[i, j], jp] * Q[j, jp]$ 
    return v[i, j] - β * jnp.sum(v[σ[i, j], :] * Q[j, :])

L_1 = jax.vmap(_L_σ, in_axes=(None, None, None, None, None, 0))
L_σ_vmap = jax.vmap(L_1, in_axes=(None, None, None, None, 0, None))

def L_σ(v, σ, params, sizes, arrays):
    w_size, y_size = sizes
    w_indices, y_indices = jnp.arange(w_size), jnp.arange(y_size)
    return L_σ_vmap(v, σ, params, arrays, w_indices, y_indices)

L_σ = jax.jit(L_σ, static_argnums=(3,))

```

Now we can define a function to compute  $v_\sigma$

```

def get_value(σ, params, sizes, arrays):
    "Get the value  $v_\sigma$  of policy  $\sigma$  by inverting the linear map  $L_\sigma$ ."

    # Unpack
    β, R, γ = params
    w_size, y_size = sizes
    w_grid, y_grid, Q = arrays

    r_σ = compute_r_σ(σ, params, sizes, arrays)

    # Reduce  $L_\sigma$  to a function in  $v$ 
    partial_L_σ = lambda v: L_σ(v, σ, params, sizes, arrays)

    return jax.scipy.sparse.linalg.bicgstab(partial_L_σ, r_σ)[0]

get_value = jax.jit(get_value, static_argnums=(2,))

```

## 12.3 Iteration

We use successive approximation for VFI.

```

def successive_approx_jax(T,                                # Operator (callable)
                          x_0,                               # Initial condition
                          tol=1e-6,                          # Error tolerance
                          max_iter=10_000):                  # Max iteration bound

    def update(inputs):
        k, x, error = inputs
        x_new = T(x)
        error = jnp.max(jnp.abs(x_new - x))

```

(continues on next page)

(continued from previous page)

```

    return k + 1, x_new, error

def condition_function(inputs):
    k, x, error = inputs
    return jnp.logical_and(error > tol, k < max_iter)

k, x, error = jax.lax.while_loop(condition_function,
                                update,
                                (1, x_0, tol + 1))

return x

successive_approx_jax = jax.jit(successive_approx_jax, static_argnums=(0,))

```

For OPI we'll add a compiled routine that computes  $T_\sigma^m v$ .

```

def iterate_policy_operator( $\sigma$ , v, m, params, sizes, arrays):

    def update(i, v):
        v = T_ $\sigma$ (v,  $\sigma$ , params, sizes, arrays)
        return v

    v = jax.lax.fori_loop(0, m, update, v)
    return v

iterate_policy_operator = jax.jit(iterate_policy_operator,
                                static_argnums=(4,))

```

## 12.4 Solvers

Now we define the solvers, which implement VFI, HPI and OPI.

Here's VFI.

```

def value_function_iteration(model, tol=1e-5):
    """
    Implements value function iteration.
    """
    params, sizes, arrays = model
    vz = jnp.zeros(sizes)
    _T = lambda v: T(v, params, sizes, arrays)
    v_star = successive_approx_jax(_T, vz, tol=tol)
    return get_greedy(v_star, params, sizes, arrays)

```

For OPI we will use a compiled JAX `lax.while_loop` operation to speed execution.

```

def opi_loop(params, sizes, arrays, m, tol, max_iter):
    """
    Implements optimistic policy iteration (see dp.quantecon.org) with
    step size m.

    """
    v_init = jnp.zeros(sizes)

```

(continues on next page)

(continued from previous page)

```

def condition_function(inputs):
    i, v, error = inputs
    return jnp.logical_and(error > tol, i < max_iter)

def update(inputs):
    i, v, error = inputs
    last_v = v
     $\sigma$  = get_greedy(v, params, sizes, arrays)
    v = iterate_policy_operator( $\sigma$ , v, m, params, sizes, arrays)
    error = jnp.max(jnp.abs(v - last_v))
    i += 1
    return i, v, error

num_iter, v, error = jax.lax.while_loop(condition_function,
                                       update,
                                       (0, v_init, tol + 1))

return get_greedy(v, params, sizes, arrays)

opi_loop = jax.jit(opi_loop, static_argnums=(1,))

```

Here's a friendly interface to OPI

```

def optimistic_policy_iteration(model, m=10, tol=1e-5, max_iter=10_000):
    params, sizes, arrays = model
     $\sigma_{\text{star}}$  = opi_loop(params, sizes, arrays, m, tol, max_iter)
    return  $\sigma_{\text{star}}$ 

```

Here's HPI.

```

def howard_policy_iteration(model, maxiter=250):
    """
    Implements Howard policy iteration (see dp.quantecon.org)
    """
    params, sizes, arrays = model
     $\sigma$  = jnp.zeros(sizes, dtype=int)
    i, error = 0, 1.0
    while error > 0 and i < maxiter:
        v_ $\sigma$  = get_value( $\sigma$ , params, sizes, arrays)
         $\sigma_{\text{new}}$  = get_greedy(v_ $\sigma$ , params, sizes, arrays)
        error = jnp.max(jnp.abs( $\sigma_{\text{new}}$  -  $\sigma$ ))
         $\sigma$  =  $\sigma_{\text{new}}$ 
        i = i + 1
        print(f"Concluded loop {i} with error {error}.")
    return  $\sigma$ 

```

## 12.5 Tests

Let's create a model for consumption, and plot the resulting optimal policy function using all the three algorithms and also check the time taken by each solver.

```
model = create_consumption_model()
# Unpack
params, sizes, arrays = model
β, R, γ = params
w_size, y_size = sizes
w_grid, y_grid, Q = arrays
```

```
2024-08-12 04:03:29.139437: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```

```
print("Starting HPI.")
start = time()
σ_star_hpi = howard_policy_iteration(model).block_until_ready()
hpi_with_compile = time() - start
print(f"HPI completed in {hpi_with_compile} seconds.")
```

Starting HPI.

```
Concluded loop 1 with error 77.
Concluded loop 2 with error 53.
Concluded loop 3 with error 28.
Concluded loop 4 with error 17.
Concluded loop 5 with error 8.
Concluded loop 6 with error 4.
Concluded loop 7 with error 1.
Concluded loop 8 with error 1.
Concluded loop 9 with error 1.
Concluded loop 10 with error 0.
HPI completed in 0.9050633907318115 seconds.
```

We run it again to get rid of compile time.

```
start = time()
σ_star_hpi = howard_policy_iteration(model).block_until_ready()
hpi_without_compile = time() - start
print(f"HPI completed in {hpi_without_compile} seconds.")
```

```
Concluded loop 1 with error 77.
Concluded loop 2 with error 53.
Concluded loop 3 with error 28.
Concluded loop 4 with error 17.
Concluded loop 5 with error 8.
Concluded loop 6 with error 4.
Concluded loop 7 with error 1.
```

```

Concluded loop 8 with error 1.
Concluded loop 9 with error 1.
Concluded loop 10 with error 0.
HPI completed in 0.10404467582702637 seconds.

```

```

fig, ax = plt.subplots()
ax.plot(w_grid, w_grid, "k--", label="45")
ax.plot(w_grid, w_grid[σ_star_hpi[:, 1]], label="$\\sigma^*_{HPI}(\cdot, y_1)$")
ax.plot(w_grid, w_grid[σ_star_hpi[:, -1]], label="$\\sigma^*_{HPI}(\cdot, y_N)$")
ax.legend()
plt.show()

```

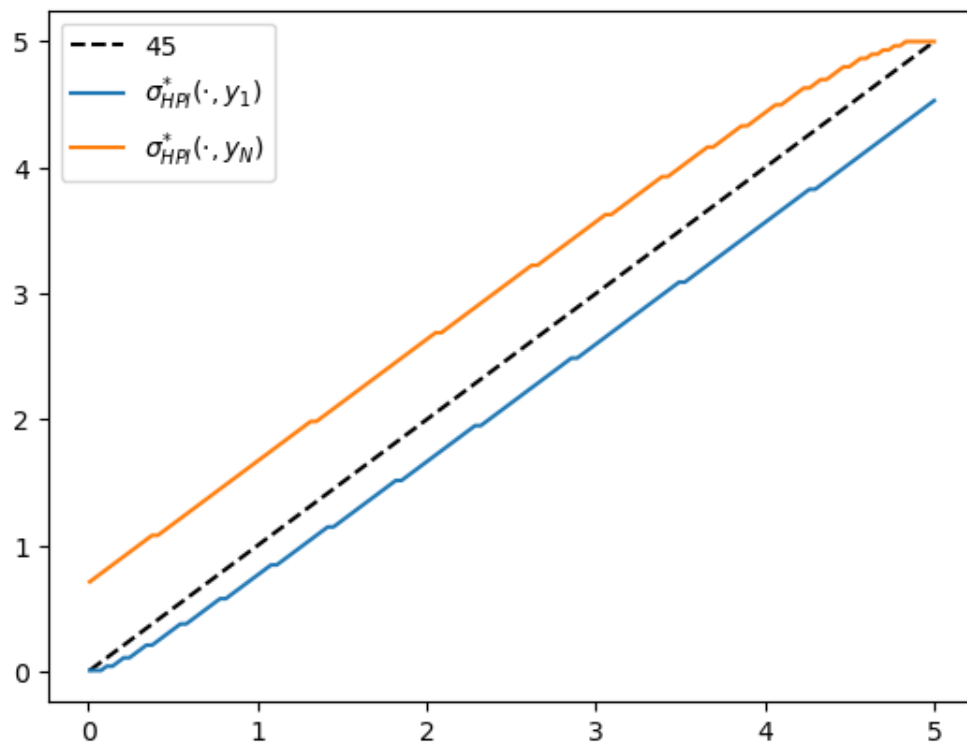


Fig. 12.1: Optimal policy function (HPI)

```

print("Starting VFI.")
start = time()
σ_star_vfi = value_function_iteration(model).block_until_ready()
vfi_with_compile = time() - start
print(f"VFI completed in {vfi_with_compile} seconds.")

```

```
Starting VFI.
```

```
VFI completed in 0.4275484085083008 seconds.
```

We run it again to eliminate compile time.



```
start = time()
sigma_star_vfi = value_function_iteration(model).block_until_ready()
vfi_without_compile = time() - start
print(f"VFI completed in {vfi_without_compile} seconds.")
```

VFI completed in 0.33431053161621094 seconds.

```
fig, ax = plt.subplots()
ax.plot(w_grid, w_grid, "k--", label="45")
ax.plot(w_grid, w_grid[sigma_star_vfi[:, 1]], label="$\\sigma^{*}_{VFI}(\cdot, y_1)$")
ax.plot(w_grid, w_grid[sigma_star_vfi[:, -1]], label="$\\sigma^{*}_{VFI}(\cdot, y_N)$")
ax.legend()
plt.show()
```

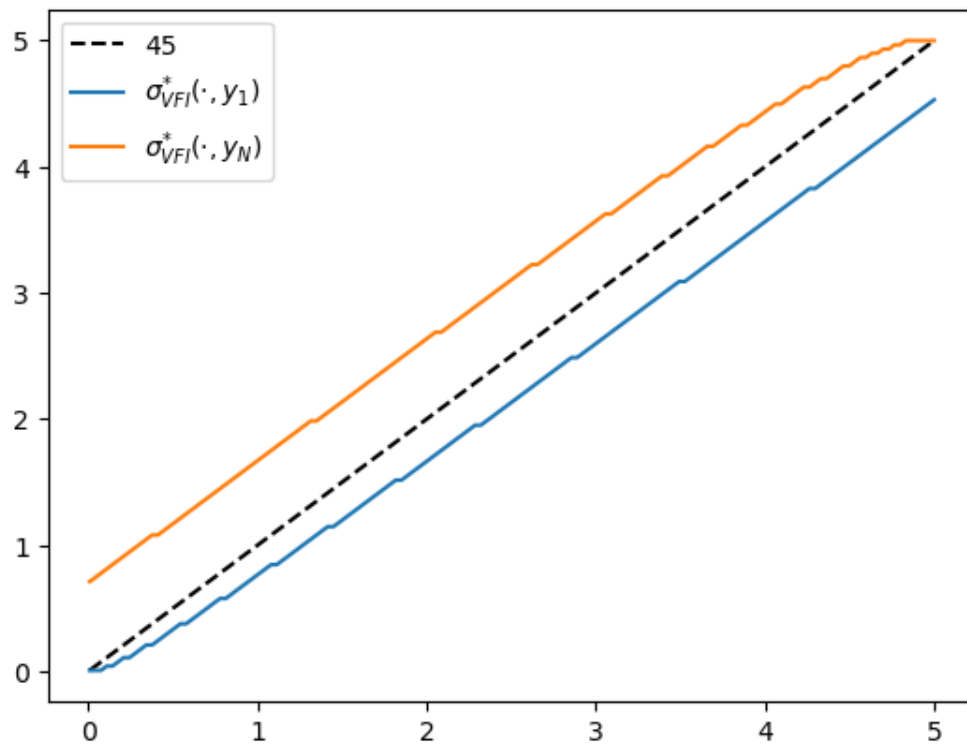


Fig. 12.2: Optimal policy function (VFI)

```
print("Starting OPI.")
start = time()
sigma_star_opi = optimistic_policy_iteration(model, m=100).block_until_ready()
opi_with_compile = time() - start
print(f"OPI completed in {opi_with_compile} seconds.")
```

Starting OPI.

```
OPI completed in 0.5606446266174316 seconds.
```

Let's run it again to get rid of compile time.

```
start = time()
sigma_star_opi = optimistic_policy_iteration(model, m=100).block_until_ready()
opi_without_compile = time() - start
print(f'OPI completed in {opi_without_compile} seconds.')
```

```
OPI completed in 0.13369250297546387 seconds.
```

```
fig, ax = plt.subplots()
ax.plot(w_grid, w_grid, "k--", label="45")
ax.plot(w_grid, w_grid[sigma_star_opi[:, 1]], label="$\\sigma^{*}_{OPI}(\\cdot, y_1)$")
ax.plot(w_grid, w_grid[sigma_star_opi[:, -1]], label="$\\sigma^{*}_{OPI}(\\cdot, y_N)$")
ax.legend()
plt.show()
```

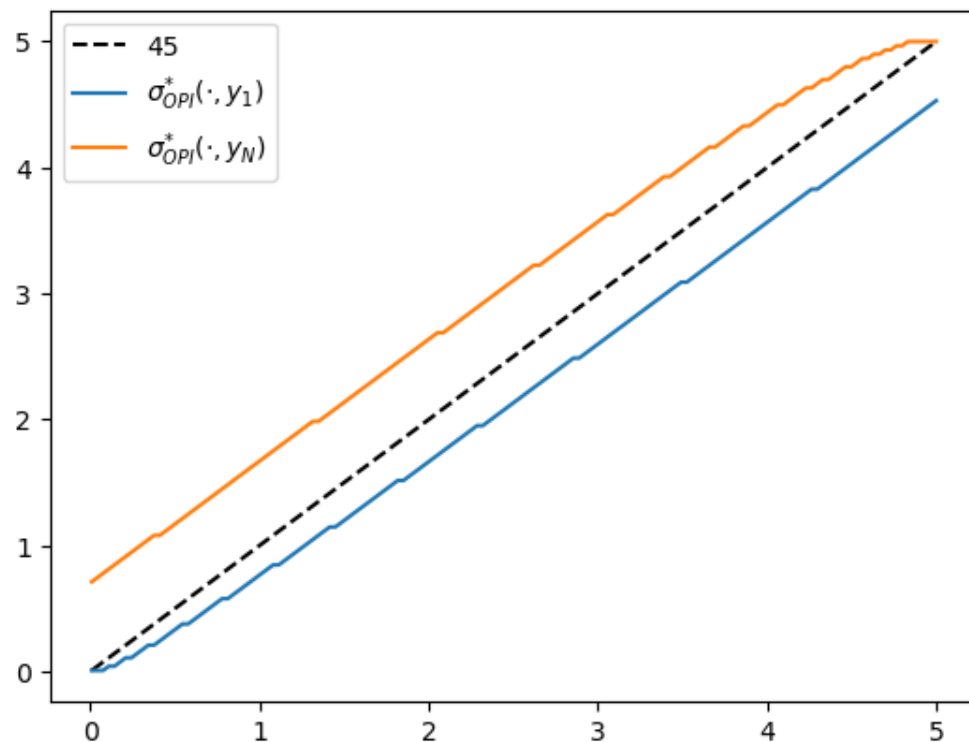


Fig. 12.3: Optimal policy function (OPI)

We observe that all the solvers produce the same output from the above three plots.

Now, let's create a plot to visualize the time differences among these algorithms.

```
def run_algorithm(algorithm, model, **kwargs):
    result = algorithm(model, **kwargs)

    # Now time it without compile time
    start = time()
    result = algorithm(model, **kwargs).block_until_ready()
    algorithm_without_compile = time() - start
    print(f"{algorithm.__name__} completed in {algorithm_without_compile:.2f} seconds.")
    ↪ return result, algorithm_without_compile
```

```
σ_pi, pi_time = run_algorithm(howard_policy_iteration, model)
σ_vfi, vfi_time = run_algorithm(value_function_iteration, model, tol=1e-5)
```

```
Concluded loop 1 with error 77.
Concluded loop 2 with error 53.
Concluded loop 3 with error 28.
Concluded loop 4 with error 17.
Concluded loop 5 with error 8.
Concluded loop 6 with error 4.
```

```
Concluded loop 7 with error 1.
Concluded loop 8 with error 1.
Concluded loop 9 with error 1.
Concluded loop 10 with error 0.
Concluded loop 1 with error 77.
Concluded loop 2 with error 53.
Concluded loop 3 with error 28.
Concluded loop 4 with error 17.
Concluded loop 5 with error 8.
Concluded loop 6 with error 4.
Concluded loop 7 with error 1.
Concluded loop 8 with error 1.
Concluded loop 9 with error 1.
Concluded loop 10 with error 0.
howard_policy_iteration completed in 0.08 seconds.
```

```
value_function_iteration completed in 0.33 seconds.
```

```
m_vals = range(5, 600, 40)
opi_times = []
for m in m_vals:
    σ_opi, opi_time = run_algorithm(optimistic_policy_iteration,
                                   model, m=m, tol=1e-5)
    opi_times.append(opi_time)
```

```
optimistic_policy_iteration completed in 0.38 seconds.
```

```
optimistic_policy_iteration completed in 0.12 seconds.
```

```
optimistic_policy_iteration completed in 0.13 seconds.
```

```
optimistic_policy_iteration completed in 0.15 seconds.
```

```
optimistic_policy_iteration completed in 0.19 seconds.
```

```
optimistic_policy_iteration completed in 0.22 seconds.
```

```
optimistic_policy_iteration completed in 0.26 seconds.
```

```
optimistic_policy_iteration completed in 0.30 seconds.
```

```
optimistic_policy_iteration completed in 0.34 seconds.
```

```
optimistic_policy_iteration completed in 0.38 seconds.
```

```
optimistic_policy_iteration completed in 0.41 seconds.
```

```
optimistic_policy_iteration completed in 0.45 seconds.
```

```
optimistic_policy_iteration completed in 0.49 seconds.
```

```
optimistic_policy_iteration completed in 0.53 seconds.
```

```
optimistic_policy_iteration completed in 0.57 seconds.
```

```
fig, ax = plt.subplots()
ax.plot(m_vals,
        jnp.full(len(m_vals), hpi_without_compile),
        lw=2, label="Howard policy iteration")
ax.plot(m_vals,
        jnp.full(len(m_vals), vfi_without_compile),
        lw=2, label="value function iteration")
ax.plot(m_vals, opi_times,
        lw=2, label="optimistic policy iteration")
ax.legend(frameon=False)
ax.set_xlabel("$m$")
ax.set_ylabel("time")
plt.show()
```

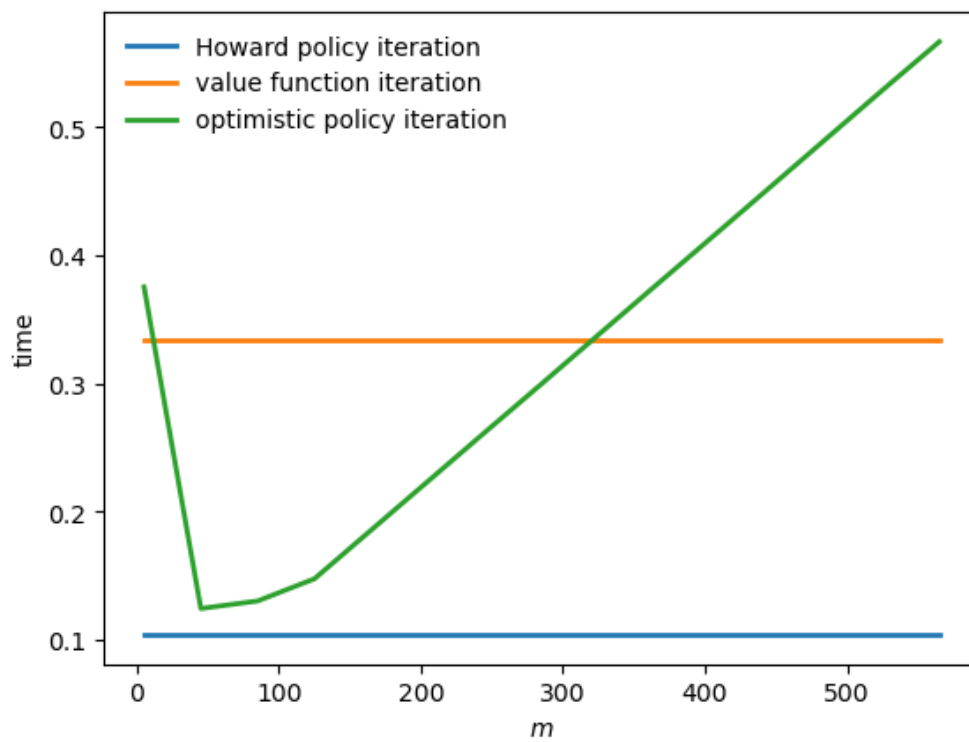


Fig. 12.4: Solver times



## SHORTEST PATHS

## GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install **Google JAX**.

## 13.1 Overview

This lecture is the extended version of the [shortest path lecture](#) using JAX. Please see that lecture for all background and notation.

Let's start by importing the libraries.

```
import numpy as np
import jax.numpy as jnp
import jax
```

Let's check the GPU we are running

```
!nvidia-smi
```

```

Mon Aug 12 04:05:49 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.1
| 5                                     |
+-----+
| GPU Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+=====+=====+=====+=====+
|    0  Tesla T4                       On  | 00000001:00:00.0 Off |
| Off |

```

(continues on next page)

(continued from previous page)

```

| N/A    63C    P0          37W /  70W |          2MiB / 16384MiB |          0%      ↵
↵Default |
|
↵ N/A |
+-----+
|
↵-----+
| Processes:
↵
| GPU  GI  CI          PID  Type  Process name          GPU ↵
↵Memory |
|      ID  ID
↵      |
|=====|
| No running processes found
↵
+-----+
↵-----+

```

## 13.2 Solving for Minimum Cost-to-Go

Let  $J(v)$  denote the minimum cost-to-go from node  $v$ , understood as the total cost from  $v$  if we take the best route.

Let's look at an algorithm for computing  $J$  and then think about how to implement it.

### 13.2.1 The Algorithm

The standard algorithm for finding  $J$  is to start an initial guess and then iterate.

This is a standard approach to solving nonlinear equations, often called the method of **successive approximations**.

Our initial guess will be

$$J_0(v) = 0 \text{ for all } v \quad (13.1)$$

Now

1. Set  $n = 0$
2. Set  $J_{n+1}(v) = \min_{w \in F_v} \{c(v, w) + J_n(w)\}$  for all  $v$
3. If  $J_{n+1}$  and  $J_n$  are not equal then increment  $n$ , go to 2

This sequence converges to  $J$ .

Let's start by defining the **distance matrix**  $Q$ .

```

inf = jnp.inf
Q = jnp.array([[inf, 1, 5, 3, inf, inf, inf],
               [inf, inf, inf, 9, 6, inf, inf],
               [inf, inf, inf, inf, inf, 2, inf],
               [inf, inf, inf, inf, inf, 4, 8],

```

(continues on next page)



(continued from previous page)

```
[inf, inf, inf, inf, inf, inf, 4],
[inf, inf, inf, inf, inf, inf, 1],
[inf, inf, inf, inf, inf, inf, 0]])
```

```
2024-08-12 04:05:50.095293: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```

Notice that the cost of staying still (on the principle diagonal) is set to

- `jnp.inf` for non-destination nodes — moving on is required.
- 0 for the destination node — here is where we stop.

Let's try with this example using python while loop and some jax vectorized code:

```
%%time

num_nodes = Q.shape[0]
J = jnp.zeros(num_nodes)

max_iter = 500
i = 0

while i < max_iter:
    next_J = jnp.min(Q + J, axis=1)
    if jnp.allclose(next_J, J):
        break
    else:
        J = next_J.copy()
        i += 1

print("The cost-to-go function is", J)
```

```
The cost-to-go function is [ 8. 10.  3.  5.  4.  1.  0.]
CPU times: user 113 ms, sys: 48.2 ms, total: 161 ms
Wall time: 193 ms
```

We can further optimize the above code by using `jax.lax.while_loop`. The extra acceleration is due to the fact that the entire operation can be optimized by the JAX compiler and launched as a single kernel on the GPU.

```
max_iter = 500
num_nodes = Q.shape[0]
J = jnp.zeros(num_nodes)
```

```
def body_fun(values):
    # Define the body function of while loop
    i, J, break_cond = values

    # Update J and break condition
    next_J = jnp.min(Q + J, axis=1)
```

(continues on next page)

(continued from previous page)

```
break_condition = jnp.allclose(next_J, J)

# Return next iteration values
return i + 1, next_J, break_condition
```

```
def cond_fun(values):
    i, J, break_condition = values
    return ~break_condition & (i < max_iter)
```

Let's see the timing for JIT compilation of the functions and runtime results.

```
%%time
jax.lax.while_loop(cond_fun, body_fun, init_val=(0, J, False))[1].block_until_ready()
```

```
CPU times: user 50.7 ms, sys: 13 ms, total: 63.8 ms
Wall time: 72.8 ms
```

```
Array([ 8., 10.,  3.,  5.,  4.,  1.,  0.], dtype=float32)
```

Now, this runs faster once we have the JIT compiled JAX version of the functions.

```
%%time
jax.lax.while_loop(cond_fun, body_fun, init_val=(0, J, False))[1].block_until_ready()
```

```
CPU times: user 681 µs, sys: 326 µs, total: 1.01 ms
Wall time: 586 µs
```

```
Array([ 8., 10.,  3.,  5.,  4.,  1.,  0.], dtype=float32)
```

---

**Note:** Large speed gains while using `jax.lax.while_loop` won't be realized unless the shortest path problem is relatively large.

---

## 13.3 Exercises

---

### Exercise 13.3.1

The text below describes a weighted directed graph.

The line `node0, node1 0.04, node8 11.11, node14 72.21` means that from `node0` we can go to

- `node1` at cost 0.04
- `node8` at cost 11.11
- `node14` at cost 72.21

No other nodes can be reached directly from `node0`.

Other lines have a similar interpretation.

Your task is to use the algorithm given above to find the optimal path and its cost.

```
%%file graph.txt
node0, node1 0.04, node8 11.11, node14 72.21
node1, node46 1247.25, node6 20.59, node13 64.94
node2, node66 54.18, node31 166.80, node45 1561.45
node3, node20 133.65, node6 2.06, node11 42.43
node4, node75 3706.67, node5 0.73, node7 1.02
node5, node45 1382.97, node7 3.33, node11 34.54
node6, node31 63.17, node9 0.72, node10 13.10
node7, node50 478.14, node9 3.15, node10 5.85
node8, node69 577.91, node11 7.45, node12 3.18
node9, node70 2454.28, node13 4.42, node20 16.53
node10, node89 5352.79, node12 1.87, node16 25.16
node11, node94 4961.32, node18 37.55, node20 65.08
node12, node84 3914.62, node24 34.32, node28 170.04
node13, node60 2135.95, node38 236.33, node40 475.33
node14, node67 1878.96, node16 2.70, node24 38.65
node15, node91 3597.11, node17 1.01, node18 2.57
node16, node36 392.92, node19 3.49, node38 278.71
node17, node76 783.29, node22 24.78, node23 26.45
node18, node91 3363.17, node23 16.23, node28 55.84
node19, node26 20.09, node20 0.24, node28 70.54
node20, node98 3523.33, node24 9.81, node33 145.80
node21, node56 626.04, node28 36.65, node31 27.06
node22, node72 1447.22, node39 136.32, node40 124.22
node23, node52 336.73, node26 2.66, node33 22.37
node24, node66 875.19, node26 1.80, node28 14.25
node25, node70 1343.63, node32 36.58, node35 45.55
node26, node47 135.78, node27 0.01, node42 122.00
node27, node65 480.55, node35 48.10, node43 246.24
node28, node82 2538.18, node34 21.79, node36 15.52
node29, node64 635.52, node32 4.22, node33 12.61
node30, node98 2616.03, node33 5.61, node35 13.95
node31, node98 3350.98, node36 20.44, node44 125.88
node32, node97 2613.92, node34 3.33, node35 1.46
node33, node81 1854.73, node41 3.23, node47 111.54
node34, node73 1075.38, node42 51.52, node48 129.45
node35, node52 17.57, node41 2.09, node50 78.81
node36, node71 1171.60, node54 101.08, node57 260.46
node37, node75 269.97, node38 0.36, node46 80.49
node38, node93 2767.85, node40 1.79, node42 8.78
node39, node50 39.88, node40 0.95, node41 1.34
node40, node75 548.68, node47 28.57, node54 53.46
node41, node53 18.23, node46 0.28, node54 162.24
node42, node59 141.86, node47 10.08, node72 437.49
node43, node98 2984.83, node54 95.06, node60 116.23
node44, node91 807.39, node46 1.56, node47 2.14
node45, node58 79.93, node47 3.68, node49 15.51
node46, node52 22.68, node57 27.50, node67 65.48
node47, node50 2.82, node56 49.31, node61 172.64
node48, node99 2564.12, node59 34.52, node60 66.44
node49, node78 53.79, node50 0.51, node56 10.89
node50, node85 251.76, node53 1.38, node55 20.10
node51, node98 2110.67, node59 23.67, node60 73.79
node52, node94 1471.80, node64 102.41, node66 123.03
node53, node72 22.85, node56 4.33, node67 88.35
node54, node88 967.59, node59 24.30, node73 238.61
```

(continues on next page)

(continued from previous page)

```

node55, node84 86.09, node57 2.13, node64 60.80
node56, node76 197.03, node57 0.02, node61 11.06
node57, node86 701.09, node58 0.46, node60 7.01
node58, node83 556.70, node64 29.85, node65 34.32
node59, node90 820.66, node60 0.72, node71 0.67
node60, node76 48.03, node65 4.76, node67 1.63
node61, node98 1057.59, node63 0.95, node64 4.88
node62, node91 132.23, node64 2.94, node76 38.43
node63, node66 4.43, node72 70.08, node75 56.34
node64, node80 47.73, node65 0.30, node76 11.98
node65, node94 594.93, node66 0.64, node73 33.23
node66, node98 395.63, node68 2.66, node73 37.53
node67, node82 153.53, node68 0.09, node70 0.98
node68, node94 232.10, node70 3.35, node71 1.66
node69, node99 247.80, node70 0.06, node73 8.99
node70, node76 27.18, node72 1.50, node73 8.37
node71, node89 104.50, node74 8.86, node91 284.64
node72, node76 15.32, node84 102.77, node92 133.06
node73, node83 52.22, node76 1.40, node90 243.00
node74, node81 1.07, node76 0.52, node78 8.08
node75, node92 68.53, node76 0.81, node77 1.19
node76, node85 13.18, node77 0.45, node78 2.36
node77, node80 8.94, node78 0.98, node86 64.32
node78, node98 355.90, node81 2.59
node79, node81 0.09, node85 1.45, node91 22.35
node80, node92 121.87, node88 28.78, node98 264.34
node81, node94 99.78, node89 39.52, node92 99.89
node82, node91 47.44, node88 28.05, node93 11.99
node83, node94 114.95, node86 8.75, node88 5.78
node84, node89 19.14, node94 30.41, node98 121.05
node85, node97 94.51, node87 2.66, node89 4.90
node86, node97 85.09
node87, node88 0.21, node91 11.14, node92 21.23
node88, node93 1.31, node91 6.83, node98 6.12
node89, node97 36.97, node99 82.12
node90, node96 23.53, node94 10.47, node99 50.99
node91, node97 22.17
node92, node96 10.83, node97 11.24, node99 34.68
node93, node94 0.19, node97 6.71, node99 32.77
node94, node98 5.91, node96 2.03
node95, node98 6.17, node99 0.27
node96, node98 3.32, node97 0.43, node99 5.87
node97, node98 0.30
node98, node99 0.33
node99,

```

Overwriting graph.txt

**Solution to Exercise 13.3.1**

First let's write a function that reads in the graph data above and builds a distance matrix.

```
num_nodes = 100
```

(continues on next page)

(continued from previous page)

```

destination_node = 99
def map_graph_to_distance_matrix(in_file):

    # First let's set of the distance matrix Q with inf everywhere
    Q = np.full((num_nodes, num_nodes), np.inf)

    # Now we read in the data and modify Q
    with open(in_file) as infile:
        for line in infile:
            elements = line.split(',')
            node = elements.pop(0)
            node = int(node[4:]) # convert node description to integer
            if node != destination_node:
                for element in elements:
                    destination, cost = element.split()
                    destination = int(destination[4:])
                    Q[node, destination] = float(cost)
            Q[destination_node, destination_node] = 0
    return jnp.array(Q)

```

Let's write a function `compute_cost_to_go` that returns  $J$  given any valid  $Q$ .

```

@jax.jit
def compute_cost_to_go(Q):
    num_nodes = Q.shape[0]
    J = jnp.zeros(num_nodes) # Initial guess
    max_iter = 500
    i = 0

    def body_fun(values):
        # Define the body function of while loop
        i, J, break_cond = values

        # Update J and break condition
        next_J = jnp.min(Q + J, axis=1)
        break_condition = jnp.allclose(next_J, J)

        # Return next iteration values
        return i + 1, next_J, break_condition

    def cond_fun(values):
        i, J, break_condition = values
        return ~break_condition & (i < max_iter)

    return jax.lax.while_loop(cond_fun, body_fun,
                             init_val=(0, J, False))[1]

```

Finally, here's a function that uses the `cost-to-go` function to obtain the optimal path (and its cost).

```

def print_best_path(J, Q):
    sum_costs = 0
    current_node = 0
    while current_node != destination_node:
        print(current_node)
        # Move to the next node and increment costs

```

(continues on next page)

(continued from previous page)

```

    next_node = jnp.argmin(Q[current_node, :] + J)
    sum_costs += Q[current_node, next_node]
    current_node = next_node
print(destination_node)
print('Cost: ', sum_costs)

```

Okay, now we have the necessary functions, let's call them to do the job we were assigned.

```
Q = map_graph_to_distance_matrix('graph.txt')
```

Let's see the timings for jitting the function and runtime results.

```
%%time
J = compute_cost_to_go(Q).block_until_ready()
```

```

CPU times: user 62.7 ms, sys: 11 ms, total: 73.6 ms
Wall time: 88.1 ms

```

Let's run again to eliminate compile time.

```
%%time
J = compute_cost_to_go(Q).block_until_ready()
```

```

CPU times: user 949 µs, sys: 68 µs, total: 1.02 ms
Wall time: 784 µs

```

```
print_best_path(J, Q)
```

```
0
```

```

8
11
18
23
33
41
53
56
57
60
67
70
73
76
85
87
88
93
94
96
97

```

(continues on next page)

(continued from previous page)

```
98  
99  
Cost: 160.55
```

The total cost of the path should agree with  $J[0]$  so let's check this.

```
J[0].item()
```

```
160.5500030517578
```





## OPTIMAL INVESTMENT

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

In addition to JAX and Anaconda, this lecture will need the following libraries:

```
!pip install --upgrade quantecon
```

We study a monopolist who faces inverse demand curve

$$P_t = a_0 - a_1 Y_t + Z_t,$$

where

- $P_t$  is price,
- $Y_t$  is output and
- $Z_t$  is a demand shock.

We assume that  $Z_t$  is a discretized AR(1) process, specified below.

Current profits are

$$P_t Y_t - c Y_t - \gamma (Y_{t+1} - Y_t)^2$$

Combining with the demand curve and writing  $y, y'$  for  $Y_t, Y_{t+1}$ , this becomes

$$r(y, z, y') := (a_0 - a_1 y + z - c)y - \gamma (y' - y)^2$$

The firm maximizes present value of expected discounted profits. The Bellman equation is

$$v(y, z) = \max_{y'} \left\{ r(y, z, y') + \beta \sum_{z'} v(y', z') Q(z, z') \right\}.$$

We discretize  $y$  to a finite grid `y_grid`.

In essence, the firm tries to choose output close to the monopolist profit maximizer, given  $Z_t$ , but is constrained by adjustment costs.

Let's begin with the following imports

```
import quantecon as qe
import jax
import jax.numpy as jnp
import matplotlib.pyplot as plt
from time import time
```

Let's check the GPU we are running

```
!nvidia-smi
```

```
Mon Aug 12 04:02:33 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5 |
+-----+
| GPU Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+=====+
|  0  Tesla T4                          On | 00000001:00:00.0 Off |
| Off |
| N/A   43C    P8              9W / 70W |      2MiB / 16384MiB |      0%    |
| Default |
|
| N/A |
+-----+

+-----+
| Processes:
|
| GPU   GI    CI        PID    Type    Process name                        GPU
| Memory |
|      ID    ID
|
|
+=====+
| No running processes found
|
+-----+
```

We will use 64 bit floats with JAX in order to increase the precision.

```
jax.config.update("jax_enable_x64", True)
```

Let's define a function to create an investment model using the given parameters.

```
def create_investment_model(
    r=0.01,                                     # Interest rate
```

(continues on next page)

(continued from previous page)

```

a_0=10.0, a_1=1.0,          # Demand parameters
y=25.0, c=1.0,              # Adjustment and unit cost
y_min=0.0, y_max=20.0, y_size=100, # Grid for output
p=0.9, v=1.0,               # AR(1) parameters
z_size=150):                # Grid size for shock
"""
A function that takes in parameters and returns an instance of Model that
contains data for the investment problem.
"""
β = 1 / (1 + r)
y_grid = jnp.linspace(y_min, y_max, y_size)
mc = qe.tauchen(z_size, p, v)
z_grid, Q = mc.state_values, mc.P

# Break up parameters into static and nonstatic components
constants = β, a_0, a_1, y, c
sizes = y_size, z_size
arrays = y_grid, z_grid, Q

# Shift arrays to the device (e.g., GPU)
arrays = tuple(map(jax.device_put, arrays))
return constants, sizes, arrays

```

Let's re-write the vectorized version of the right-hand side of the Bellman equation (before maximization), which is a 3D array representing

$$B(y, z, y') = r(y, z, y') + \beta \sum_{z'} v(y', z') Q(z, z')$$

for all  $(y, z, y')$ .

```

def B(v, constants, sizes, arrays):
    """
    A vectorized version of the right-hand side of the Bellman equation
    (before maximization)
    """

    # Unpack
    β, a_0, a_1, y, c = constants
    y_size, z_size = sizes
    y_grid, z_grid, Q = arrays

    # Compute current rewards r(y, z, yp) as array r[i, j, ip]
    y = jnp.reshape(y_grid, (y_size, 1, 1)) # y[i] -> y[i, j, ip]
    z = jnp.reshape(z_grid, (1, z_size, 1)) # z[j] -> z[i, j, ip]
    yp = jnp.reshape(y_grid, (1, 1, y_size)) # yp[ip] -> yp[i, j, ip]
    r = (a_0 - a_1 * y + z - c) * y - y * (yp - y)**2

    # Calculate continuation rewards at all combinations of (y, z, yp)
    v = jnp.reshape(v, (1, 1, y_size, z_size)) # v[ip, jp] -> v[i, j, ip, jp]
    Q = jnp.reshape(Q, (1, z_size, 1, z_size)) # Q[j, jp] -> Q[i, j, ip, jp]
    EV = jnp.sum(v * Q, axis=3) # sum over last index jp

    # Compute the right-hand side of the Bellman equation
    return r + β * EV

```

(continues on next page)

(continued from previous page)

```
# Create a jitted function
B = jax.jit(B, static_argnums=(2,))
```

We define a function to compute the current rewards  $r_\sigma$  given policy  $\sigma$ , which is defined as the vector

$$r_\sigma(y, z) := r(y, z, \sigma(y, z))$$

```
def compute_r_sigma(sigma, constants, sizes, arrays):
    """
    Compute the array r_sigma[i, j] = r[i, j, sigma[i, j]], which gives current
    rewards given policy sigma.
    """

    # Unpack model
    beta, a_0, a_1, y, c = constants
    y_size, z_size = sizes
    y_grid, z_grid, Q = arrays

    # Compute r_sigma[i, j]
    y = jnp.reshape(y_grid, (y_size, 1))
    z = jnp.reshape(z_grid, (1, z_size))
    yp = y_grid[sigma]
    r_sigma = (a_0 - a_1 * y + z - c) * y - y * (yp - y)**2

    return r_sigma

# Create the jitted function
compute_r_sigma = jax.jit(compute_r_sigma, static_argnums=(2,))
```

Define the Bellman operator.

```
def T(v, constants, sizes, arrays):
    """The Bellman operator."""
    return jnp.max(B(v, constants, sizes, arrays), axis=2)

T = jax.jit(T, static_argnums=(2,))
```

The following function computes a v-greedy policy.

```
def get_greedy(v, constants, sizes, arrays):
    "Computes a v-greedy policy, returned as a set of indices."
    return jnp.argmax(B(v, constants, sizes, arrays), axis=2)

get_greedy = jax.jit(get_greedy, static_argnums=(2,))
```

Define the  $\sigma$ -policy operator.

```
def T_sigma(v, sigma, constants, sizes, arrays):
    """The sigma-policy operator."""

    # Unpack model
    beta, a_0, a_1, y, c = constants
    y_size, z_size = sizes
```

(continues on next page)

(continued from previous page)

```

y_grid, z_grid, Q = arrays

r_σ = compute_r_σ(σ, constants, sizes, arrays)

# Compute the array v[σ[i, j], jp]
zp_idx = jnp.arange(z_size)
zp_idx = jnp.reshape(zp_idx, (1, 1, z_size))
σ = jnp.reshape(σ, (y_size, z_size, 1))
V = v[σ, zp_idx]

# Convert Q[j, jp] to Q[i, j, jp]
Q = jnp.reshape(Q, (1, z_size, z_size))

# Calculate the expected sum E_jp v[σ[i, j], jp] * Q[i, j, jp]
Ev = jnp.sum(V * Q, axis=2)

return r_σ + β * Ev

T_σ = jax.jit(T_σ, static_argnums=(3,))

```

Next, we want to compute the lifetime value of following policy  $\sigma$ .

This lifetime value is a function  $v_\sigma$  that satisfies

$$v_\sigma(y, z) = r_\sigma(y, z) + \beta \sum_{z'} v_\sigma(\sigma(y, z), z') Q(z, z')$$

We wish to solve this equation for  $v_\sigma$ .

Suppose we define the linear operator  $L_\sigma$  by

$$(L_\sigma v)(y, z) = v(y, z) - \beta \sum_{z'} v(\sigma(y, z), z') Q(z, z')$$

With this notation, the problem is to solve for  $v$  via

$$(L_\sigma v)(y, z) = r_\sigma(y, z)$$

In vector form this is  $L_\sigma v = r_\sigma$ , which tells us that the function we seek is

$$v_\sigma = L_\sigma^{-1} r_\sigma$$

JAX allows us to solve linear systems defined in terms of operators; the first step is to define the function  $L_\sigma$ .

```

def L_σ(v, σ, constants, sizes, arrays):

    β, a_0, a_1, γ, c = constants
    y_size, z_size = sizes
    y_grid, z_grid, Q = arrays

    # Set up the array v[σ[i, j], jp]
    zp_idx = jnp.arange(z_size)
    zp_idx = jnp.reshape(zp_idx, (1, 1, z_size))
    σ = jnp.reshape(σ, (y_size, z_size, 1))
    V = v[σ, zp_idx]

    # Expand Q[j, jp] to Q[i, j, jp]

```

(continues on next page)

(continued from previous page)

```

Q = jnp.reshape(Q, (1, z_size, z_size))

# Compute and return  $v[i, j] - \beta \sum_{jp} v[\sigma[i, j], jp] * Q[j, jp]$ 
return v -  $\beta$  * jnp.sum(V * Q, axis=2)

L_σ = jax.jit(L_σ, static_argnums=(3,))

```

Now we can define a function to compute  $v_\sigma$

```

def get_value(σ, constants, sizes, arrays):

    # Unpack
    β, a_0, a_1, γ, c = constants
    y_size, z_size = sizes
    y_grid, z_grid, Q = arrays

    r_σ = compute_r_σ(σ, constants, sizes, arrays)

    # Reduce  $L_\sigma$  to a function in  $v$ 
    partial_L_σ = lambda v: L_σ(v, σ, constants, sizes, arrays)

    return jax.scipy.sparse.linalg.bicgstab(partial_L_σ, r_σ)[0]

get_value = jax.jit(get_value, static_argnums=(2,))

```

We use successive approximation for VFI.

```

def successive_approx_jax(T,                # Operator (callable)
                        x_0,                # Initial condition
                        tol=1e-6,           # Error tolerance
                        max_iter=10_000):    # Max iteration bound

    def body_fun(k_x_err):
        k, x, error = k_x_err
        x_new = T(x)
        error = jnp.max(jnp.abs(x_new - x))
        return k + 1, x_new, error

    def cond_fun(k_x_err):
        k, x, error = k_x_err
        return jnp.logical_and(error > tol, k < max_iter)

    k, x, error = jax.lax.while_loop(cond_fun, body_fun, (1, x_0, tol + 1))
    return x

successive_approx_jax = jax.jit(successive_approx_jax, static_argnums=(0,))

```

For OPI we'll add a compiled routine that computes  $T_\sigma^m v$ .

```

def iterate_policy_operator(σ, v, m, params, sizes, arrays):

    def update(i, v):
        v = T_σ(v, σ, params, sizes, arrays)
        return v

    v = jax.lax.fori_loop(0, m, update, v)

```

(continues on next page)

(continued from previous page)

```

    return v

iterate_policy_operator = jax.jit(iterate_policy_operator,
                                  static_argnums=(4,))

```

Finally, we introduce the solvers that implement VFI, HPI and OPI.

```

def value_function_iteration(model, tol=1e-5):
    """
    Implements value function iteration.
    """
    params, sizes, arrays = model
    vz = jnp.zeros(sizes)
    _T = lambda v: T(v, params, sizes, arrays)
    v_star = successive_approx_jax(_T, vz, tol=tol)
    return get_greedy(v_star, params, sizes, arrays)

```

For OPI we will use a compiled JAX `lax.while_loop` operation to speed execution.

```

def opi_loop(params, sizes, arrays, m, tol, max_iter):
    """
    Implements optimistic policy iteration (see dp.quantecon.org) with
    step size m.
    """
    v_init = jnp.zeros(sizes)

    def condition_function(inputs):
        i, v, error = inputs
        return jnp.logical_and(error > tol, i < max_iter)

    def update(inputs):
        i, v, error = inputs
        last_v = v
        sigma = get_greedy(v, params, sizes, arrays)
        v = iterate_policy_operator(sigma, v, m, params, sizes, arrays)
        error = jnp.max(jnp.abs(v - last_v))
        i += 1
        return i, v, error

    num_iter, v, error = jax.lax.while_loop(condition_function,
                                           update,
                                           (0, v_init, tol + 1))

    return get_greedy(v, params, sizes, arrays)

opi_loop = jax.jit(opi_loop, static_argnums=(1,))

```

Here's a friendly interface to OPI

```

def optimistic_policy_iteration(model, m=10, tol=1e-5, max_iter=10_000):
    params, sizes, arrays = model
    sigma_star = opi_loop(params, sizes, arrays, m, tol, max_iter)
    return sigma_star

```

Here's HPI

```
def howard_policy_iteration(model, maxiter=250):
    """
    Implements Howard policy iteration (see dp.quantecon.org)
    """
    params, sizes, arrays = model
     $\sigma$  = jnp.zeros(sizes, dtype=int)
    i, error = 0, 1.0
    while error > 0 and i < maxiter:
        v_ $\sigma$  = get_value( $\sigma$ , params, sizes, arrays)
         $\sigma_{\text{new}}$  = get_greedy(v_ $\sigma$ , params, sizes, arrays)
        error = jnp.max(jnp.abs( $\sigma_{\text{new}}$  -  $\sigma$ ))
         $\sigma$  =  $\sigma_{\text{new}}$ 
        i = i + 1
        print(f"Concluded loop {i} with error {error}.")
    return  $\sigma$ 
```

```
model = create_investment_model()
constants, sizes, arrays = model
 $\beta$ , a_0, a_1, y, c = constants
y_size, z_size = sizes
y_grid, z_grid, Q = arrays
```

```
print("Starting HPI.")
%time  $\sigma_{\text{star\_hpi}}$  = howard_policy_iteration(model).block_until_ready()
```

```
# Now time it without compile time
start = time()
 $\sigma_{\text{star\_hpi}}$  = howard_policy_iteration(model).block_until_ready()
hpi_without_compile = time() - start
print( $\sigma_{\text{star\_hpi}}$ )
print(f"HPI completed in {hpi_without_compile} seconds.")
```

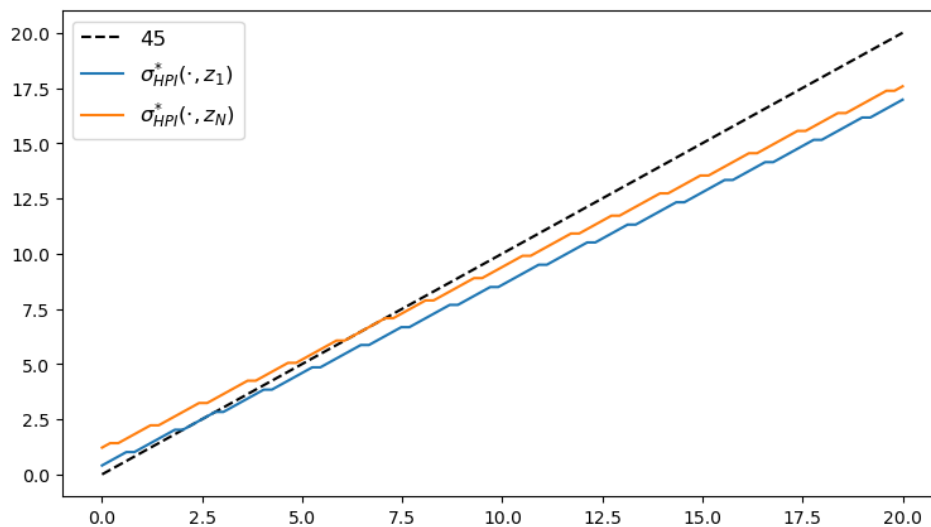
```
Concluded loop 1 with error 50.
Concluded loop 2 with error 26.
Concluded loop 3 with error 17.
Concluded loop 4 with error 10.
Concluded loop 5 with error 7.
Concluded loop 6 with error 4.
```

```
Concluded loop 7 with error 3.
Concluded loop 8 with error 1.
Concluded loop 9 with error 1.
Concluded loop 10 with error 1.
Concluded loop 11 with error 1.
Concluded loop 12 with error 0.
[[ 2  2  2 ...  6  6  6]
 [ 3  3  3 ...  7  7  7]
 [ 4  4  4 ...  7  7  7]
 ...
 [82 82 82 ... 86 86 86]
 [83 83 83 ... 86 86 86]
 [84 84 84 ... 87 87 87]]
HPI completed in 0.11728167533874512 seconds.
```



Here's the plot of the Howard policy, as a function of  $y$  at the highest and lowest values of  $z$ .

```
fig, ax = plt.subplots(figsize=(9, 5))
ax.plot(y_grid, y_grid, "k--", label="45")
ax.plot(y_grid, y_grid[σ_star_hpi[:, 1]], label="$\\sigma^{*}_{HPI}(\cdot, z_1)$")
ax.plot(y_grid, y_grid[σ_star_hpi[:, -1]], label="$\\sigma^{*}_{HPI}(\cdot, z_N)$")
ax.legend(fontsize=12)
plt.show()
```



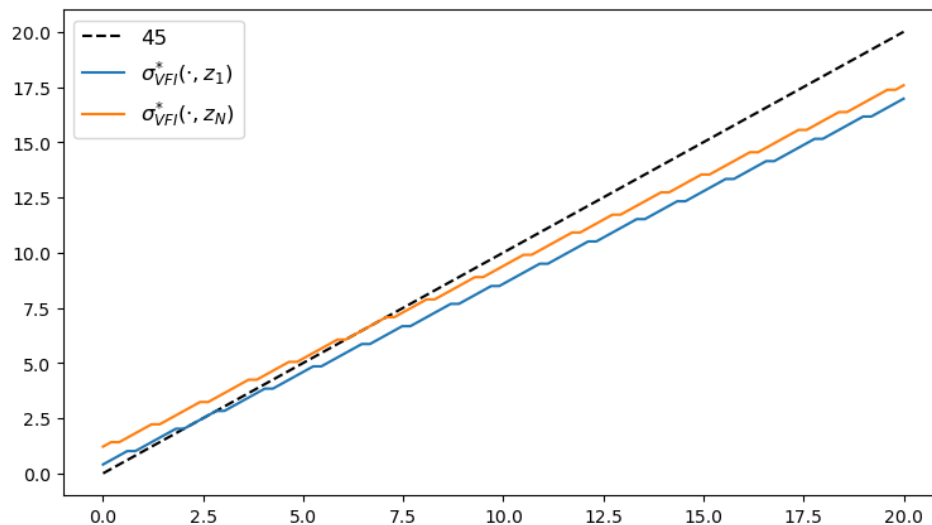
```
print("Starting VFI.")
%time σ_star_vfi = value_function_iteration(model).block_until_ready()
```

```
# Now time it without compile time
start = time()
σ_star_vfi = value_function_iteration(model).block_until_ready()
vfi_without_compile = time() - start
print(σ_star_vfi)
print(f"VFI completed in {vfi_without_compile} seconds.")
```

```
[[ 2  2  2 ...  6  6  6]
 [ 3  3  3 ...  7  7  7]
 [ 4  4  4 ...  7  7  7]
 ...
 [82 82 82 ... 86 86 86]
 [83 83 83 ... 86 86 86]
 [84 84 84 ... 87 87 87]]
VFI completed in 0.43529224395751953 seconds.
```

Here's the plot of the VFI, as a function of  $y$  at the highest and lowest values of  $z$ .

```
fig, ax = plt.subplots(figsize=(9, 5))
ax.plot(y_grid, y_grid, "k--", label="45")
ax.plot(y_grid, y_grid[σ_star_vfi[:, 1]], label="$\\sigma^{*}_{VFI}(\cdot, z_1)$")
ax.plot(y_grid, y_grid[σ_star_vfi[:, -1]], label="$\\sigma^{*}_{VFI}(\cdot, z_N)$")
ax.legend(fontsize=12)
plt.show()
```



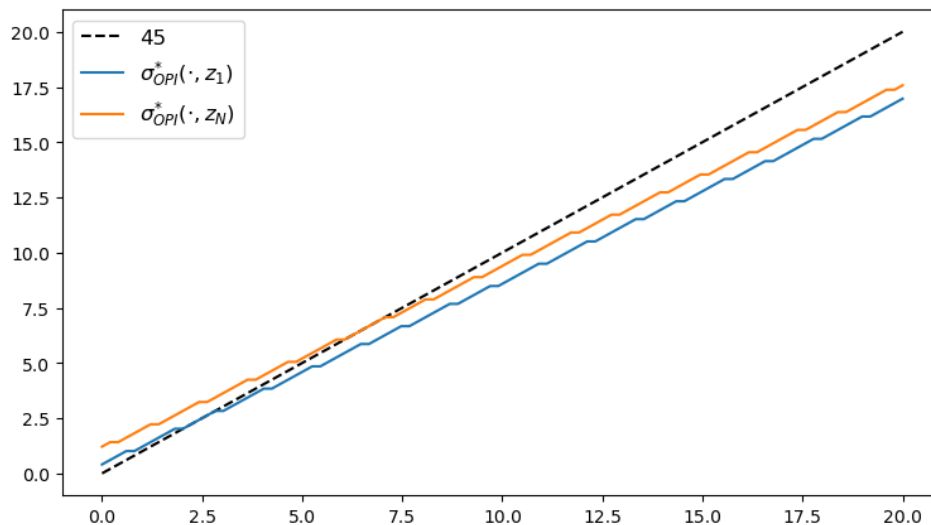
```
print("Starting OPI.")
%time σ_star_opi = optimistic_policy_iteration(model, m=100).block_until_ready()
```

```
# Now time it without compile time
start = time()
σ_star_opi = optimistic_policy_iteration(model, m=100).block_until_ready()
opi_without_compile = time() - start
print(σ_star_opi)
print(f"OPI completed in {opi_without_compile} seconds.")
```

```
[[ 2  2  2 ...  6  6  6]
 [ 3  3  3 ...  7  7  7]
 [ 4  4  4 ...  7  7  7]
 ...
 [82 82 82 ... 86 86 86]
 [83 83 83 ... 86 86 86]
 [84 84 84 ... 87 87 87]]
OPI completed in 0.2242722511291504 seconds.
```

Here's the plot of the optimal policy, as a function of  $y$  at the highest and lowest values of  $z$ .

```
fig, ax = plt.subplots(figsize=(9, 5))
ax.plot(y_grid, y_grid, "k--", label="45")
ax.plot(y_grid, y_grid[σ_star_opi[:, 1]], label="$\\sigma^{*}_{OPI}(\\cdot, z_1)$")
ax.plot(y_grid, y_grid[σ_star_opi[:, -1]], label="$\\sigma^{*}_{OPI}(\\cdot, z_N)$")
ax.legend(fontsize=12)
plt.show()
```



We observe that all the solvers produce the same output from the above three plots.

Let's plot the time taken by each of the solvers and compare them.

```
m_vals = range(5, 600, 40)
```

```
print("Running Howard policy iteration.")
%time σ_hpi = howard_policy_iteration(model).block_until_ready()
```

```
Running Howard policy iteration.
Concluded loop 1 with error 50.
Concluded loop 2 with error 26.
Concluded loop 3 with error 17.
Concluded loop 4 with error 10.
```

```
Concluded loop 5 with error 7.
Concluded loop 6 with error 4.
Concluded loop 7 with error 3.
Concluded loop 8 with error 1.
Concluded loop 9 with error 1.
Concluded loop 10 with error 1.
Concluded loop 11 with error 1.
Concluded loop 12 with error 0.
CPU times: user 102 ms, sys: 875 μs, total: 103 ms
Wall time: 94.7 ms
```

```
# Now time it without compile time
start = time()
σ_hpi = howard_policy_iteration(model).block_until_ready()
hpi_without_compile = time() - start
print(f"HPI completed in {hpi_without_compile} seconds.")
```

```
Concluded loop 1 with error 50.
Concluded loop 2 with error 26.
Concluded loop 3 with error 17.
```

(continues on next page)

(continued from previous page)

```

Concluded loop 4 with error 10.
Concluded loop 5 with error 7.
Concluded loop 6 with error 4.
Concluded loop 7 with error 3.
Concluded loop 8 with error 1.
Concluded loop 9 with error 1.
Concluded loop 10 with error 1.
Concluded loop 11 with error 1.
Concluded loop 12 with error 0.
HPI completed in 0.09428048133850098 seconds.

```

```

print("Running value function iteration.")
%time σ_vfi = value_function_iteration(model, tol=1e-5).block_until_ready()

```

```
Running value function iteration.
```

```

CPU times: user 429 ms, sys: 11.8 ms, total: 441 ms
Wall time: 436 ms

```

```

# Now time it without compile time
start = time()
σ_vfi = value_function_iteration(model, tol=1e-5).block_until_ready()
vfi_without_compile = time() - start
print(f"VFI completed in {vfi_without_compile} seconds.")

```

```
VFI completed in 0.43514227867126465 seconds.
```

```

opi_times = []
for m in m_vals:
    print(f"Running optimistic policy iteration with m={m}.")
    σ_opi = optimistic_policy_iteration(model, m=m, tol=1e-5).block_until_ready()

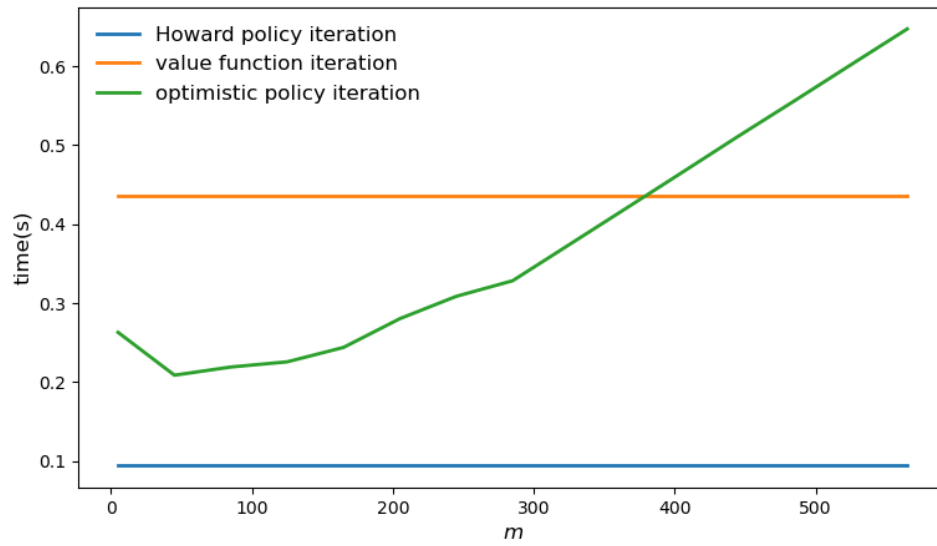
    # Now time it without compile time
    start = time()
    σ_opi = optimistic_policy_iteration(model, m=m, tol=1e-5).block_until_ready()
    opi_without_compile = time() - start
    print(f"OPI with m={m} completed in {opi_without_compile} seconds.")
    opi_times.append(opi_without_compile)

```

```

fig, ax = plt.subplots(figsize=(9, 5))
ax.plot(m_vals, jnp.full(len(m_vals), hpi_without_compile),
        lw=2, label="Howard policy iteration")
ax.plot(m_vals, jnp.full(len(m_vals), vfi_without_compile),
        lw=2, label="value function iteration")
ax.plot(m_vals, opi_times, lw=2, label="optimistic policy iteration")
ax.legend(fontsize=12, frameon=False)
ax.set_xlabel("$m$", fontsize=12)
ax.set_ylabel("time(s)", fontsize=12)
plt.show()

```





## INVENTORY MANAGEMENT MODEL

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

This lecture provides a JAX implementation of a model in [Dynamic Programming](#).

In addition to JAX and Anaconda, this lecture will need the following libraries:

```
!pip install --upgrade quantecon
```

## 15.1 A model with constant discounting

We study a firm where a manager tries to maximize shareholder value.

To simplify the problem, we assume that the firm only sells one product.

Letting  $\pi_t$  be profits at time  $t$  and  $r > 0$  be the interest rate, the value of the firm is

$$V_0 = \sum_{t \geq 0} \beta^t \pi_t \quad \text{where} \quad \beta := \frac{1}{1+r}.$$

Suppose the firm faces exogenous demand process  $(D_t)_{t \geq 0}$ .

We assume  $(D_t)_{t \geq 0}$  is IID with common distribution  $\phi \in (Z_+)$ .

Inventory  $(X_t)_{t \geq 0}$  of the product obeys

$$X_{t+1} = f(X_t, D_{t+1}, A_t) \quad \text{where} \quad f(x, a, d) := (x - d) \vee 0 + a.$$

The term  $A_t$  is units of stock ordered this period, which take one period to arrive.

We assume that the firm can store at most  $K$  items at one time.

Profits are given by

$$\pi_t := X_t \wedge D_{t+1} - cA_t - \kappa 1\{A_t > 0\}.$$

We take the minimum of current stock and demand because orders in excess of inventory are assumed to be lost rather than back-filled.

Here  $c$  is unit product cost and  $\kappa$  is a fixed cost of ordering inventory.

We can map our inventory problem into a dynamic program with state space  $X := \{0, \dots, K\}$  and action space  $A := X$ .

The feasible correspondence  $\Gamma$  is

$$\Gamma(x) := \{0, \dots, K - x\},$$

which represents the set of feasible orders when the current inventory state is  $x$ .

The reward function is expected current profits, or

$$r(x, a) := \sum_{d \geq 0} (x \wedge d) \phi(d) - ca - \kappa 1\{a > 0\}.$$

The stochastic kernel (i.e., state-transition probabilities) from the set of feasible state-action pairs is

$$P(x, a, x') := P\{f(x, a, D) = x'\} \quad \text{when } D \sim \phi.$$

When discounting is constant, the Bellman equation takes the form

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{d \geq 0} v(f(x, a, d)) \phi(d) \right\} \quad (15.1)$$

## 15.2 Time varying discount rates

We wish to consider a more sophisticated model with time-varying discounting.

This time variation accommodates non-constant interest rates.

To this end, we replace the constant  $\beta$  in (15.1) with a stochastic process  $(\beta_t)$  where

- $\beta_t = 1/(1 + r_t)$  and
- $r_t$  is the interest rate at time  $t$

We suppose that the dynamics can be expressed as  $\beta_t = \beta(Z_t)$ , where the exogenous process  $(Z_t)_{t \geq 0}$  is a Markov chain on  $Z$  with Markov matrix  $Q$ .

After relabeling inventory  $X_t$  as  $Y_t$  and  $x$  as  $y$ , the Bellman equation becomes

$$v(y, z) = \max_{a \in \Gamma(y)} B((y, z), a, v)$$

where

$$B((y, z), a, v) = r(y, a) + \beta(z) \sum_{d, z'} v(f(y, a, d), z') \phi(d) Q(z, z'). \quad (15.2)$$

We set

$$R(y, a, y') := P\{f(y, a, d) = y'\} \quad \text{when } D \sim \phi,$$

Now  $R(y, a, y')$  is the probability of realizing next period inventory level  $y'$  when the current level is  $y$  and the action is  $a$ .

Hence we can rewrite (15.2) as

$$B((y, z), a, v) = r(y, a) + \beta(z) \sum_{y', z'} v(y', z') Q(z, z') R(y, a, y').$$

Let's begin with the following imports



```

import quantecon as qe
import jax
import jax.numpy as jnp
import numpy as np
import matplotlib.pyplot as plt
from collections import namedtuple
import numba
from numba import prange
from time import time

```

Let's check the GPU we are running

```
!nvidia-smi
```

```

Mon Aug 12 03:26:54 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5 |
+-----+
| GPU   Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+-----+
|  0  Tesla T4                               On | 00000001:00:00.0 Off |
| Off |
| N/A   42C    P8              9W / 70W |  2MiB / 16384MiB |      0%
| Default |
|
| N/A |
+-----+
+-----+
| Processes:
|
| GPU   GI    CI        PID   Type   Process name                      GPU
| Memory |
| ID     ID
|
+-----+
| No running processes found
+-----+

```

We will use 64 bit floats with JAX in order to increase the precision.

```
jax.config.update("jax_enable_x64", True)
```

Let's define a model to represent the inventory management.

```
# NamedTuple Model
Model = namedtuple("Model", ("c", "κ", "p", "z_vals", "Q"))
```

We need the following successive approximation function.

```
def successive_approx(T,                # Operator (callable)
                      x_0,              # Initial condition
                      tolerance=1e-6,   # Error tolerance
                      max_iter=10_000,  # Max iteration bound
                      print_step=25,    # Print at multiples
                      verbose=False):
    x = x_0
    error = tolerance + 1
    k = 1
    while error > tolerance and k <= max_iter:
        x_new = T(x)
        error = jnp.max(jnp.abs(x_new - x))
        if verbose and k % print_step == 0:
            print(f"Completed iteration {k} with error {error}.")
        x = x_new
        k += 1
    if error > tolerance:
        print(f"Warning: Iteration hit upper bound {max_iter}.")
    elif verbose:
        print(f"Terminated successfully in {k} iterations.")
    return x
```

```
@jax.jit
def demand_pdf(p, d):
    return (1 - p)**d * p
```

```
K = 100
D_MAX = 101
```

Let's define a function to create an inventory model using the given parameters.

```
def create_sdd_inventory_model(
    p=0.98, v=0.002, n_z=100, b=0.97,      # Z state parameters
    c=0.2, κ=0.8, p=0.6,                  # firm and demand parameters
    use_jax=True):
    mc = qe.tauchen(n_z, p, v)
    z_vals, Q = mc.state_values + b, mc.P
    if use_jax:
        z_vals, Q = map(jnp.array, (z_vals, Q))
    return Model(c=c, κ=κ, p=p, z_vals=z_vals, Q=Q)
```

Here's the function B on the right-hand side of the Bellman equation.

```
@jax.jit
def B(x, i_z, a, v, model):
    """
    The function  $B(x, z, a, v) = r(x, a) + \beta(z) \sum_{x'} v(x') P(x, a, x')$ .
    """
    c, κ, p, z_vals, Q = model
```

(continues on next page)

(continued from previous page)

```

z = z_vals[i_z]
d_vals = jnp.arange(D_MAX)
phi_vals = demand_pdf(p, d_vals)
revenue = jnp.sum(jnp.minimum(x, d_vals)*phi_vals)
profit = revenue - c * a - kappa * (a > 0)
v_R = jnp.sum(v[jnp.maximum(x - d_vals, 0) + a].T * phi_vals, axis=1)
cv = jnp.sum(v_R*Q[i_z])
return profit + z * cv

```

We need to vectorize this function so that we can use it efficiently in JAX.

We apply a sequence of vmap operations to vectorize appropriately in each argument.

```
B_vec_a = jax.vmap(B, in_axes=(None, None, 0, None, None))
```

```

@jax.jit
def B2(x, i_z, v, model):
    """
    The function  $B(x, z, a, v) = r(x, a) + \beta(z) \sum_{x'} v(x') P(x, a, x')$ .
    """
    c, kappa, p, z_vals, Q = model
    a_vals = jnp.arange(K)
    res = B_vec_a(x, i_z, a_vals, v, model)
    return jnp.where(a_vals < K - x + 1, res, -jnp.inf)

```

```

B2_vec_z = jax.vmap(B2, in_axes=(None, 0, None, None))
B2_vec_z_x = jax.vmap(B2_vec_z, in_axes=(0, None, None, None))

```

Next we define the Bellman operator.

```

@jax.jit
def T(v, model):
    """The Bellman operator."""
    c, kappa, p, z_vals, Q = model
    i_z_range = jnp.arange(len(z_vals))
    x_range = jnp.arange(K + 1)
    res = B2_vec_z_x(x_range, i_z_range, v, model)
    return jnp.max(res, axis=2)

```

The following function computes a v-greedy policy.

```

@jax.jit
def get_greedy(v, model):
    """Get a v-greedy policy. Returns a zero-based array."""
    c, kappa, p, z_vals, Q = model
    i_z_range = jnp.arange(len(z_vals))
    x_range = jnp.arange(K + 1)
    res = B2_vec_z_x(x_range, i_z_range, v, model)
    return jnp.argmax(res, axis=2)

```

Here's code to solve the model using value function iteration.

```

def solve_inventory_model(v_init, model):
    """Use successive_approx to get v_star and then compute greedy."""

```

(continues on next page)

(continued from previous page)

```

v_star = successive_approx(lambda v: T(v, model), v_init, verbose=True)
σ_star = get_greedy(v_star, model)
return v_star, σ_star

```

Now let's create an instance and solve it.

```

model = create_sdd_inventory_model()
c, κ, p, z_vals, Q = model
n_z = len(z_vals)
v_init = jnp.zeros((K + 1, n_z), dtype=float)

```

```

2024-08-12 03:26:55.300737: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.

```

```

start = time()
v_star, σ_star = solve_inventory_model(v_init, model)
jax_time_with_compile = time() - start
print("Jax compile plus execution time = ", jax_time_with_compile)

```

```
Completed iteration 25 with error 0.5613828428334688.
```

```
Completed iteration 50 with error 0.37764643476880266.
```

```
Completed iteration 75 with error 0.2272706235969011.
```

```
Completed iteration 100 with error 0.12872204940709508.
```

```
Completed iteration 125 with error 0.06744149371262154.
```

```
Completed iteration 150 with error 0.03037463954767361.
```

```
Completed iteration 175 with error 0.01423099032950148.
```

```
Completed iteration 200 with error 0.007396776219316337.
```

```
Completed iteration 225 with error 0.0039122383045793185.
```

```
Completed iteration 250 with error 0.002068091416653317.
```

```
Completed iteration 275 with error 0.001092307533355097.
```

Completed iteration 300 with error 0.0005766427105911021.

Completed iteration 325 with error 0.00030433217072101115.

Completed iteration 350 with error 0.00016059073674767887.

Completed iteration 375 with error 8.473334524694565e-05.

Completed iteration 400 with error 4.4706045166265085e-05.

Completed iteration 425 with error 2.3586619946058818e-05.

Completed iteration 450 with error 1.2443945934137446e-05.

Completed iteration 475 with error 6.5651783245357365e-06.

Completed iteration 500 with error 3.463639430378862e-06.

Completed iteration 525 with error 1.827332347659194e-06.

Terminated successfully in 550 iterations.  
Jax compile plus execution time = 5.884882926940918

Let's run again to get rid of the compile time.

```
start = time()
v_star, sigma_star = solve_inventory_model(v_init, model)
jax_time_without_compile = time() - start
print("Jax execution time = ", jax_time_without_compile)
```

Completed iteration 25 with error 0.5613828428334688.

Completed iteration 50 with error 0.37764643476880266.

Completed iteration 75 with error 0.2272706235969011.

Completed iteration 100 with error 0.12872204940709508.

Completed iteration 125 with error 0.06744149371262154.

Completed iteration 150 with error 0.03037463954767361.

Completed iteration 175 with error 0.01423099032950148.

```
Completed iteration 200 with error 0.007396776219316337.
```

```
Completed iteration 225 with error 0.0039122383045793185.
```

```
Completed iteration 250 with error 0.002068091416653317.
```

```
Completed iteration 275 with error 0.001092307533355097.
```

```
Completed iteration 300 with error 0.0005766427105911021.
```

```
Completed iteration 325 with error 0.00030433217072101115.
```

```
Completed iteration 350 with error 0.00016059073674767887.
```

```
Completed iteration 375 with error 8.473334524694565e-05.
```

```
Completed iteration 400 with error 4.4706045166265085e-05.
```

```
Completed iteration 425 with error 2.3586619946058818e-05.
```

```
Completed iteration 450 with error 1.2443945934137446e-05.
```

```
Completed iteration 475 with error 6.5651783245357365e-06.
```

```
Completed iteration 500 with error 3.463639430378862e-06.
```

```
Completed iteration 525 with error 1.827332347659194e-06.
```

```
Terminated successfully in 550 iterations.  
Jax execution time = 5.286510229110718
```

```
z_mc = ge.MarkovChain(Q, z_vals)
```

```
def sim_inventories(ts_length, X_init=0):  
    """Simulate given the optimal policy."""  
    global p, z_mc  
    i_z = z_mc.simulate_indices(ts_length, init=1)  
    X = np.zeros(ts_length, dtype=np.int32)  
    X[0] = X_init  
    rand = np.random.default_rng().geometric(p=p, size=ts_length-1) - 1  
    for t in range(ts_length-1):  
        X[t+1] = np.maximum(X[t] - rand[t], 0) +  $\sigma_{\text{star}}$ [X[t], i_z[t]]  
    return X, z_vals[i_z]
```

```
def plot_ts(ts_length=400, fontsize=10):
    X, Z = sim_inventories(ts_length)
    fig, axes = plt.subplots(2, 1, figsize=(9, 5.5))

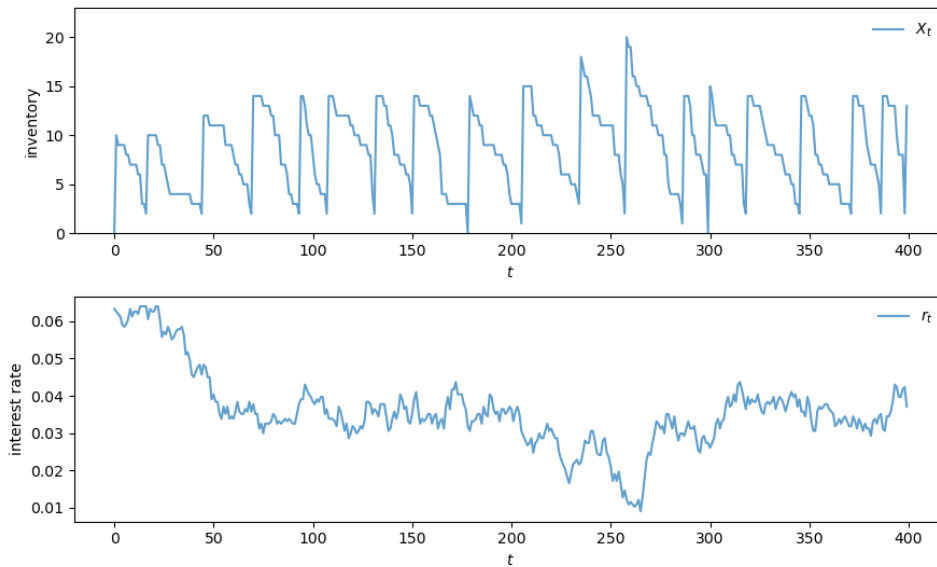
    ax = axes[0]
    ax.plot(X, label=r"$X_t$", alpha=0.7)
    ax.set_xlabel(r"$t$", fontsize=fontsize)
    ax.set_ylabel("inventory", fontsize=fontsize)
    ax.legend(fontsize=fontsize, frameon=False)
    ax.set_ylim(0, np.max(X)+3)

    # calculate interest rate from discount factors
    r = (1 / Z) - 1

    ax = axes[1]
    ax.plot(r, label=r"$r_t$", alpha=0.7)
    ax.set_xlabel(r"$t$", fontsize=fontsize)
    ax.set_ylabel("interest rate", fontsize=fontsize)
    ax.legend(fontsize=fontsize, frameon=False)

    plt.tight_layout()
    plt.show()
```

```
plot_ts()
```



## 15.3 Numba implementation

Let's try the same operations in Numba in order to compare the speed.

```
@numba.njit
def demand_pdf_numba(p, d):
    return (1 - p)**d * p

@numba.njit
def B_numba(x, i_z, a, v, model):
    """
    The function  $B(x, z, a, v) = r(x, a) + \beta(z) \sum_{x'} v(x') P(x, a, x')$ .
    """
    c, k, p, z_vals, Q = model
    z = z_vals[i_z]
    d_vals = np.arange(D_MAX)
    phi_vals = demand_pdf_numba(p, d_vals)
    revenue = np.sum(np.minimum(x, d_vals) * phi_vals)
    profit = revenue - c * a - k * (a > 0)
    v_R = np.sum(v[np.maximum(x - d_vals, 0) + a].T * phi_vals, axis=1)
    cv = np.sum(v_R * Q[i_z])
    return profit + z * cv

@numba.njit(parallel=True)
def T_numba(v, model):
    """The Bellman operator."""
    c, k, p, z_vals, Q = model
    new_v = np.empty_like(v)
    for i_z in prange(len(z_vals)):
        for x in prange(K+1):
            v_1 = np.array([B_numba(x, i_z, a, v, model)
                           for a in range(K-x+1)])
            new_v[x, i_z] = np.max(v_1)
    return new_v

@numba.njit(parallel=True)
def get_greedy_numba(v, model):
    """Get a v-greedy policy. Returns a zero-based array."""
    c, k, p, z_vals, Q = model
    n_z = len(z_vals)
    sigma_star = np.zeros((K+1, n_z), dtype=np.int32)
    for i_z in prange(n_z):
        for x in range(K+1):
            v_1 = np.array([B_numba(x, i_z, a, v, model)
                           for a in range(K-x+1)])
            sigma_star[x, i_z] = np.argmax(v_1)
    return sigma_star

def solve_inventory_model_numba(v_init, model):
    """Use successive_approx to get v_star and then compute greedy."""
    v_star = successive_approx(lambda v: T_numba(v, model), v_init, verbose=True)
    sigma_star = get_greedy_numba(v_star, model)
    return v_star, sigma_star
```



```

model = create_sdd_inventory_model(use_jax=False)
c, k, p, z_vals, Q = model
n_z = len(z_vals)
v_init = np.zeros((K + 1, n_z), dtype=float)

```

```

start = time()
v_star_numba, sigma_star_numba = solve_inventory_model_numba(v_init, model)
numba_time_with_compile = time() - start
print("Numba compile plus execution time = ", numba_time_with_compile)

```

Completed iteration 25 with error 0.5613828428334706.

Completed iteration 50 with error 0.37764643476879556.

Completed iteration 75 with error 0.22727062359689398.

Completed iteration 100 with error 0.12872204940708798.

Completed iteration 125 with error 0.06744149371262864.

Completed iteration 150 with error 0.030374639547666504.

Completed iteration 175 with error 0.01423099032948727.

Completed iteration 200 with error 0.007396776219316337.

Completed iteration 225 with error 0.003912238304593529.

Completed iteration 250 with error 0.002068091416653317.

Completed iteration 275 with error 0.0010923075333622023.

Completed iteration 300 with error 0.0005766427105911021.

Completed iteration 325 with error 0.0003043321707281166.

Completed iteration 350 with error 0.00016059073676188973.

Completed iteration 375 with error 8.473334525405107e-05.

Completed iteration 400 with error 4.470604518047594e-05.

Completed iteration 425 with error 2.3586619960269672e-05.

Completed iteration 450 with error 1.2443945934137446e-05.

Completed iteration 475 with error 6.565178331641164e-06.

Completed iteration 500 with error 3.4636394445897167e-06.

Completed iteration 525 with error 1.827332347659194e-06.

Terminated successfully in 550 iterations.

Numba compile plus execution time = 951.5428130626678

Let's run again to eliminate the compile time.

```
start = time()
v_star_numba, sigma_star_numba = solve_inventory_model_numba(v_init, model)
numba_time_without_compile = time() - start
print("Numba execution time = ", numba_time_without_compile)
```

Completed iteration 25 with error 0.5613828428334706.

Completed iteration 50 with error 0.37764643476879556.

Completed iteration 75 with error 0.22727062359689398.

Completed iteration 100 with error 0.12872204940708798.

Completed iteration 125 with error 0.06744149371262864.

Completed iteration 150 with error 0.030374639547666504.

Completed iteration 175 with error 0.01423099032948727.

Completed iteration 200 with error 0.007396776219316337.

Completed iteration 225 with error 0.003912238304593529.

Completed iteration 250 with error 0.002068091416653317.

Completed iteration 275 with error 0.0010923075333622023.

Completed iteration 300 with error 0.0005766427105911021.

```
Completed iteration 325 with error 0.0003043321707281166.
```

```
Completed iteration 350 with error 0.00016059073676188973.
```

```
Completed iteration 375 with error 8.473334525405107e-05.
```

```
Completed iteration 400 with error 4.470604518047594e-05.
```

```
Completed iteration 425 with error 2.3586619960269672e-05.
```

```
Completed iteration 450 with error 1.2443945934137446e-05.
```

```
Completed iteration 475 with error 6.565178331641164e-06.
```

```
Completed iteration 500 with error 3.4636394445897167e-06.
```

```
Completed iteration 525 with error 1.827332347659194e-06.
```

```
Terminated successfully in 550 iterations.
```

```
Numba execution time = 949.7083163261414
```

Let's verify that the Numba and JAX implementations converge to the same solution.

```
np.allclose(v_star_numba, v_star)
```

```
True
```

Here's the speed comparison.

```
print("JAX vectorized implementation is "
      f"{numba_time_without_compile/jax_time_without_compile} faster "
      "than Numba's parallel implementation")
```

```
JAX vectorized implementation is 179.64749431420256 faster than Numba's parallel
implementation
```



## ENDOGENOUS GRID METHOD

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

## 16.1 Overview

In this lecture we use the endogenous grid method (EGM) to solve a basic income fluctuation (optimal savings) problem.

Background on the endogenous grid method can be found in an [earlier QuantEcon lecture](#).

Here we focus on providing an efficient JAX implementation.

In addition to JAX and Anaconda, this lecture will need the following libraries:

```
!pip install --upgrade quantecon
```

```
import quantecon as qe
import matplotlib.pyplot as plt
import numpy as np
import jax
import jax.numpy as jnp
import numba
from time import time
```

Let’s check the GPU we are running

```
!nvidia-smi
```

```
Mon Aug 12 03:24:14 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5                                    |
+-----+-----+-----+
| 5                                    |
```

(continues on next page)

(continued from previous page)

```

| GPU Name Persistence-M | Bus-Id Disp.A | Volatile Uncorr.
↳ ECC |
| Fan Temp Perf Pwr:Usage/Cap | Memory-Usage | GPU-Util  ↳
↳ Compute M. |
|
|                               |                               |                               |
=====+=====+=====+=====+
|  0 Tesla T4 On | 00000001:00:00.0 Off |                               |  ↳
↳ Off |
| N/A 43C P8 11W / 70W | 2MiB / 16384MiB | 0%  ↳
↳ Default |
|
|                               |                               |                               |
↳ N/A |
+-----+-----+-----+-----+
↳-----+
+-----+
↳-----+
| Processes:  ↳
↳ |
| GPU GI CI PID Type Process name GPU  ↳
↳ Memory |
| ID ID
↳ |
|                               |                               |                               |
=====+=====+=====+=====+
| No running processes found  ↳
↳ |
+-----+
↳-----+

```

We use 64 bit floating point numbers for extra precision.

```
jax.config.update("jax_enable_x64", True)
```

## 16.2 Setup

We consider a household that chooses a state-contingent consumption plan  $\{c_t\}_{t \geq 0}$  to maximize

$$\mathbb{E} \sum_{t=0}^{\infty} \beta^t u(c_t)$$

subject to

$$a_{t+1} \leq R(a_t - c_t) + Y_{t+1}, \quad c_t \geq 0, \quad a_t \geq 0 \quad t = 0, 1, \dots$$

Here  $R = 1 + r$  where  $r$  is the interest rate.

The income process  $\{Y_t\}$  is a **Markov chain** generated by stochastic matrix  $P$ .

The matrix  $P$  and the grid of values taken by  $Y_t$  are obtained by discretizing the AR(1) process

$$Y_{t+1} = \rho Y_t + \nu \epsilon_{t+1}$$

where  $\{\epsilon_t\}$  is IID and standard normal.

Utility has the CRRA specification

$$u(c) = \frac{c^{1-\gamma}}{1-\gamma}$$

The following function stores default parameter values for the income fluctuation problem and creates suitable arrays.

```
def ifp(R=1.01,          # gross interest rate
        beta=0.99,       # discount factor
        gamma=1.5,       # CRRA preference parameter
        s_max=16,        # savings grid max
        s_size=200,      # savings grid size
        rho=0.99,        # income persistence
        v=0.02,          # income volatility
        y_size=25):      # income grid size

    # require R * beta < 1 for convergence
    assert R * beta < 1, "Stability condition failed."
    # Create income Markov chain
    mc = qe.tauchen(y_size, rho, v)
    y_grid, P = jnp.exp(mc.state_values), mc.P
    # Shift to JAX arrays
    P, y_grid = jax.device_put((P, y_grid))
    s_grid = jnp.linspace(0, s_max, s_size)
    # Pack and return
    constants = beta, R, gamma
    sizes = s_size, y_size
    arrays = s_grid, y_grid, P
    return constants, sizes, arrays
```

## 16.3 Solution method

Let  $S = \mathbb{R}_+ \times Y$  be the set of possible values for the state  $(a_t, Y_t)$ .

We aim to compute an optimal consumption policy  $\sigma^*: S \rightarrow \mathbb{R}$ , under which dynamics are given by

$$c_t = \sigma^*(a_t, Y_t) \quad \text{and} \quad a_{t+1} = R(a_t - c_t) + Y_{t+1}$$

In this section we discuss how we intend to solve for this policy.

### 16.3.1 Euler equation

The Euler equation for the optimization problem is

$$u'(c_t) = \max \{ \beta R \mathbb{E}_t u'(c_{t+1}), u'(a_t) \}$$

An explanation for this expression can be found [here](#).

We rewrite the Euler equation in functional form

$$(u' \circ \sigma)(a, y) = \max \{ \beta R \mathbb{E}_y (u' \circ \sigma)[R(a - \sigma(a, y)) + \hat{Y}, \hat{Y}], u'(a) \}$$

where  $(u' \circ \sigma)(a, y) := u'(\sigma(a, y))$  and  $\sigma$  is a consumption policy.

For given consumption policy  $\sigma$ , we define  $(K\sigma)(a, y)$  as the unique  $c \in [0, a]$  that solves

$$u'(c) = \max \left\{ \beta R \mathbb{E}_y(u' \circ \sigma) [R(a - c) + \hat{Y}, \hat{Y}], u'(a) \right\} \quad (16.1)$$

It can be shown that

1. iterating with  $K$  computes an optimal policy and
2. if  $\sigma$  is increasing in its first argument, then so is  $K\sigma$

Hence below we always assume that  $\sigma$  is increasing in its first argument.

The EGM is a technique for computing the update  $K\sigma$  given  $\sigma$  along a grid of asset values.

Notice that, since  $u'(a) \rightarrow \infty$  as  $a \downarrow 0$ , the second term in the max in (16.3.1) dominates for sufficiently small  $a$ .

Also, again using (16.3.1), we have  $c = a$  for all such  $a$ .

Hence, for sufficiently small  $a$ ,

$$u'(a) \geq \beta R \mathbb{E}_y(u' \circ \sigma) [\hat{Y}, \hat{Y}]$$

Equality holds at  $\bar{a}(y)$  given by

$$\bar{a}(y) = (u')^{-1} \left\{ \beta R \mathbb{E}_y(u' \circ \sigma) [\hat{Y}, \hat{Y}] \right\}$$

We can now write

$$u'(c) = \begin{cases} \beta R \mathbb{E}_y(u' \circ \sigma) [R(a - c) + \hat{Y}, \hat{Y}] & \text{if } a > \bar{a}(y) \\ u'(a) & \text{if } a \leq \bar{a}(y) \end{cases}$$

Equivalently, we can state that the  $c$  satisfying  $c = (K\sigma)(a, y)$  obeys

$$c = \begin{cases} (u')^{-1} \left\{ \beta R \mathbb{E}_y(u' \circ \sigma) [R(a - c) + \hat{Y}, \hat{Y}] \right\} & \text{if } a > \bar{a}(y) \\ a & \text{if } a \leq \bar{a}(y) \end{cases} \quad (16.2)$$

We begin with an *exogenous* grid of saving values  $0 = s_0 < \dots < s_{N-1}$

Using the exogenous savings grid, and a fixed value of  $y$ , we create an *endogenous* asset grid  $a_0, \dots, a_{N-1}$  and a consumption grid  $c_0, \dots, c_{N-1}$  as follows.

First we set  $a_0 = c_0 = 0$ , since zero consumption is an optimal (in fact the only) choice when  $a = 0$ .

Then, for  $i > 0$ , we compute

$$c_i = (u')^{-1} \left\{ \beta R \mathbb{E}_y(u' \circ \sigma) [Rs_i + \hat{Y}, \hat{Y}] \right\} \quad \text{for all } i \quad (16.3)$$

and we set

$$a_i = s_i + c_i$$

We claim that each pair  $a_i, c_i$  obeys (16.3.2).

Indeed, since  $s_i > 0$ , choosing  $c_i$  according to (16.3.3) gives

$$c_i = (u')^{-1} \left\{ \beta R \mathbb{E}_y(u' \circ \sigma) [Rs_i + \hat{Y}, \hat{Y}] \right\} \geq \bar{a}(y)$$

where the inequality uses the fact that  $\sigma$  is increasing in its first argument.

If we now take  $a_i = s_i + c_i$  we get  $a_i > \bar{a}(y)$ , so the pair  $(a_i, c_i)$  satisfies

$$c_i = (u')^{-1} \left\{ \beta R \mathbb{E}_y(u' \circ \sigma) [R(a_i - c_i) + \hat{Y}, \hat{Y}] \right\} \quad \text{and} \quad a_i > \bar{a}(y)$$

Hence (16.3.2) holds.

We are now ready to iterate with  $K$ .



### 16.3.2 JAX version

First we define a vectorized operator  $K$  based on the EGM.

Notice in the code below that

- we avoid all loops and any mutation of arrays
- the function is pure (no globals, no mutation of inputs)

```
def K_egm(a_in, σ_in, constants, sizes, arrays):
    """
    The vectorized operator K using EGM.

    """

    # Unpack
    β, R, γ = constants
    s_size, y_size = sizes
    s_grid, y_grid, P = arrays

    def u_prime(c):
        return c**(-γ)

    def u_prime_inv(u):
        return u**(-1/γ)

    # Linearly interpolate σ(a, y)
    def σ(a, y):
        return jnp.interp(a, a_in[:, y], σ_in[:, y])
    σ_vec = jnp.vectorize(σ)

    # Broadcast and vectorize
    y_hat = jnp.reshape(y_grid, (1, 1, y_size))
    y_hat_idx = jnp.reshape(jnp.arange(y_size), (1, 1, y_size))
    s = jnp.reshape(s_grid, (s_size, 1, 1))
    P = jnp.reshape(P, (1, y_size, y_size))

    # Evaluate consumption choice
    a_next = R * s + y_hat
    σ_next = σ_vec(a_next, y_hat_idx)
    up = u_prime(σ_next)
    E = jnp.sum(up * P, axis=-1)
    c = u_prime_inv(β * R * E)

    # Set up a column vector with zero in the first row and ones elsewhere
    e_0 = jnp.ones(s_size) - jnp.identity(s_size)[:, 0]
    e_0 = jnp.reshape(e_0, (s_size, 1))

    # The policy is computed consumption with the first row set to zero
    σ_out = c * e_0

    # Compute a_out by a = s + c
    a_out = np.reshape(s_grid, (s_size, 1)) + σ_out

    return a_out, σ_out
```

Then we use `jax.jit` to compile  $K$ .

We use `static_argnums` to allow a recompile whenever sizes changes, since the compiler likes to specialize on

shapes.

```
K_egm_jax = jax.jit(K_egm, static_argnums=(3,))
```

Next we define a successive approximator that repeatedly applies  $K$ .

```
def successive_approx_jax(model,
    tol=1e-5,
    max_iter=100_000,
    verbose=True,
    print_skip=25):

    # Unpack
    constants, sizes, arrays = model
    β, R, γ = constants
    s_size, y_size = sizes
    s_grid, y_grid, P = arrays

    # Initial condition is to consume all in every state
    σ_init = jnp.repeat(s_grid, y_size)
    σ_init = jnp.reshape(σ_init, (s_size, y_size))
    a_init = jnp.copy(σ_init)
    a_vec, σ_vec = a_init, σ_init

    i = 0
    error = tol + 1

    while i < max_iter and error > tol:
        a_new, σ_new = K_egm_jax(a_vec, σ_vec, constants, sizes, arrays)
        error = jnp.max(jnp.abs(σ_vec - σ_new))
        i += 1
        if verbose and i % print_skip == 0:
            print(f"Error at iteration {i} is {error}.")
            a_vec, σ_vec = jnp.copy(a_new), jnp.copy(σ_new)

    if error > tol:
        print("Failed to converge!")
    else:
        print(f"\nConverged in {i} iterations.")

    return a_new, σ_new
```

### 16.3.3 Numba version

Below we provide a second set of code, which solves the same model with Numba.

The purpose of this code is to cross-check our results from the JAX version, as well as to do a runtime comparison.

Most readers will want to skip ahead to the next section, where we solve the model and run the cross-check.

```
@numba.jit
def K_egm_nb(a_in, σ_in, constants, sizes, arrays):
    """
    The operator K using Numba.

    """
```

(continues on next page)

(continued from previous page)

```

# Simplify names
β, R, γ = constants
s_size, y_size = sizes
s_grid, y_grid, P = arrays

def u_prime(c):
    return c**(-γ)

def u_prime_inv(u):
    return u**(-1/γ)

# Linear interpolation of policy using endogenous grid
def σ(a, z):
    return np.interp(a, a_in[:, z], σ_in[:, z])

# Allocate memory for new consumption array
σ_out = np.zeros_like(σ_in)
a_out = np.zeros_like(σ_out)

for i, s in enumerate(s_grid[1:]):
    i += 1
    for z in range(y_size):
        expect = 0.0
        for z_hat in range(y_size):
            expect += u_prime(σ(R * s + y_grid[z_hat], z_hat)) * \
                P[z, z_hat]
        c = u_prime_inv(β * R * expect)
        σ_out[i, z] = c
        a_out[i, z] = s + c

return a_out, σ_out

```

```

def successive_approx_numba(model,          # Class with model information
                           tol=1e-5,
                           max_iter=100_000,
                           verbose=True,
                           print_skip=25):

    # Unpack
    constants, sizes, arrays = model
    s_size, y_size = sizes
    # make NumPy versions of arrays
    arrays = tuple(map(np.array, arrays))
    s_grid, y_grid, P = arrays

    σ_init = np.repeat(s_grid, y_size)
    σ_init = np.reshape(σ_init, (s_size, y_size))
    a_init = np.copy(σ_init)
    a_vec, σ_vec = a_init, σ_init

    # Set up loop
    i = 0
    error = tol + 1

```

(continues on next page)

(continued from previous page)

```

while i < max_iter and error > tol:
    a_new, σ_new = K_egm_nb(a_vec, σ_vec, constants, sizes, arrays)
    error = np.max(np.abs(σ_vec - σ_new))
    i += 1
    if verbose and i % print_skip == 0:
        print(f"Error at iteration {i} is {error}.")
    a_vec, σ_vec = np.copy(a_new), np.copy(σ_new)

if error > tol:
    print("Failed to converge!")
else:
    print(f"\nConverged in {i} iterations.")

return a_new, σ_new

```

## 16.4 Solutions

Here we solve the IFP with JAX and Numba.

We will compare both the outputs and the execution time.

### 16.4.1 Outputs

```
model = ifp()
```

```

2024-08-12 03:24:14.865027: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.

```

Here's a first run of the JAX code.

```

%%time
a_star_egm_jax, σ_star_egm_jax = successive_approx_jax(model,
                                                         print_skip=1000)

```

```
Error at iteration 1000 is 6.472028596182788e-05.
```

```
Error at iteration 2000 is 1.2994575430580468e-05.
```

```
Converged in 2192 iterations.
```

```
CPU times: user 2.51 s, sys: 877 ms, total: 3.38 s
```

```
Wall time: 2.61 s
```

Next let's solve the same IFP with Numba.

```
%%time
a_star_egm_nb,  $\sigma$ _star_egm_nb = successive_approx_numba(model,
                                                             print_skip=1000)
```

```
Error at iteration 1000 is 6.472028596182788e-05.
```

```
Error at iteration 2000 is 1.2994575430802513e-05.
```

```
Converged in 2192 iterations.
CPU times: user 54.7 s, sys: 12.8 ms, total: 54.7 s
Wall time: 54.7 s
```

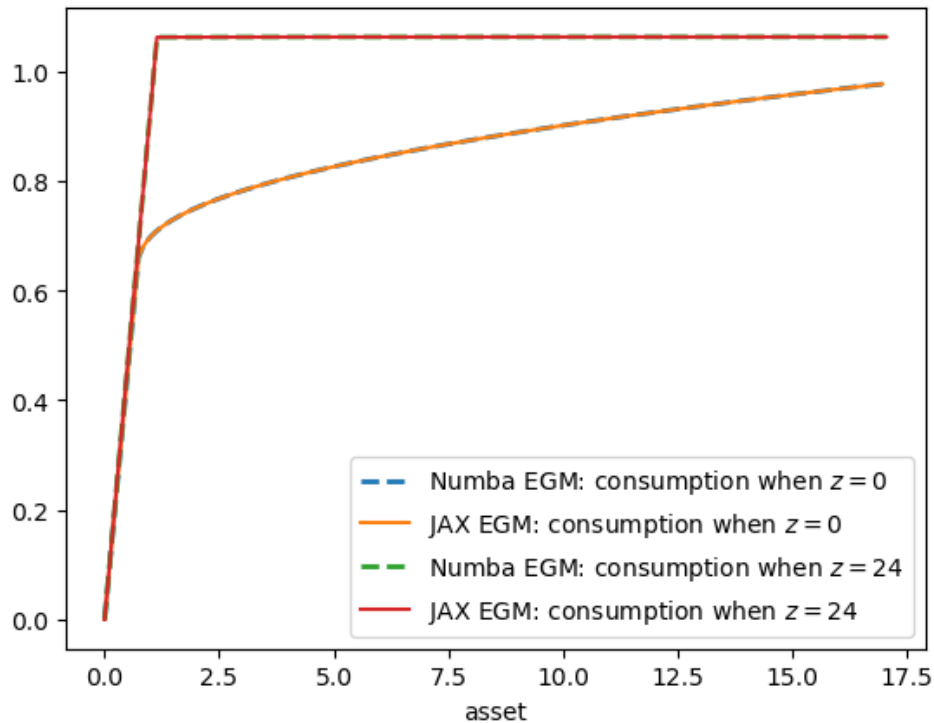
Now let's check the outputs in a plot to make sure they are the same.

```
constants, sizes, arrays = model
 $\beta$ , R,  $\gamma$  = constants
s_size, y_size = sizes
s_grid, y_grid, P = arrays

fig, ax = plt.subplots()

for z in (0, y_size-1):
    ax.plot(a_star_egm_nb[:, z],
             $\sigma$ _star_egm_nb[:, z],
            '--', lw=2,
            label=f"Numba EGM: consumption when  $\$z=\{z\}\$")
    ax.plot(a_star_egm_jax[:, z],
             $\sigma$ _star_egm_jax[:, z],
            label=f"JAX EGM: consumption when  $\$z=\{z\}\$")

ax.set_xlabel('asset')
plt.legend()
plt.show()$$ 
```



### 16.4.2 Timing

Now let's compare execution time of the two methods.

```
start = time()
a_star_egm_jax, sigma_star_egm_jax = successive_approx_jax(model,
                                                           print_skip=1000)
jax_time_without_compile = time() - start
print("Jax execution time = ", jax_time_without_compile)
```

```
Error at iteration 1000 is 6.472028596182788e-05.
```

```
Error at iteration 2000 is 1.2994575430580468e-05.
```

```
Converged in 2192 iterations.
Jax execution time = 2.0226478576660156
```

```
start = time()
a_star_egm_nb, sigma_star_egm_nb = successive_approx_numba(model,
                                                            print_skip=1000)
numba_time_without_compile = time() - start
print("Numba execution time = ", numba_time_without_compile)
```

```
Error at iteration 1000 is 6.472028596182788e-05.
```

```
Error at iteration 2000 is 1.2994575430802513e-05.
```

```
Converged in 2192 iterations.  
Numba execution time = 52.99245834350586
```

```
jax_time_without_compile / numba_time_without_compile
```

```
0.03816859834195421
```

The JAX code is significantly faster, as expected.

This difference will increase when more features (and state variables) are added to the model.





## **Part V**

# **Macroeconomic Models**



## DEFAULT RISK AND INCOME FLUCTUATIONS

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

In addition to JAX and Anaconda, this lecture will need the following libraries:

```
!pip install quantecon
```

### 17.1 Overview

This lecture computes versions of Arellano’s [[Arellano, 2008](#)] model of sovereign default.

The model describes interactions among default risk, output, and an equilibrium interest rate that includes a premium for endogenous default risk.

The decision maker is a government of a small open economy that borrows from risk-neutral foreign creditors.

The foreign lenders must be compensated for default risk.

The government borrows and lends abroad in order to smooth the consumption of its citizens.

The government repays its debt only if it wants to, but declining to pay has adverse consequences.

The interest rate on government debt adjusts in response to the state-dependent default probability chosen by government.

The model yields outcomes that help interpret sovereign default experiences, including

- countercyclical interest rates on sovereign debt
- countercyclical trade balances
- high volatility of consumption relative to output

Notably, long recessions caused by bad draws in the income process increase the government’s incentive to default.

This can lead to

- spikes in interest rates
- temporary losses of access to international credit markets

- large drops in output, consumption, and welfare
- large capital outflows during recessions

Such dynamics are consistent with experiences of many countries.

Let's start with some imports:

```
import matplotlib.pyplot as plt
import quantecon as qe
import random

import jax
import jax.numpy as jnp
from collections import namedtuple
```

Let's check the GPU we are running

```
!nvidia-smi
```

```
Mon Aug 12 03:21:35 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5 |
+-----+
| GPU   Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| Fan  Temp   Perf           Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+-----+
|  0  Tesla T4                       On      | 00000001:00:00.0 Off  |
| Off |
| N/A   39C    P8              9W / 70W | 2MiB / 16384MiB |      0%    |
| Default |
|
| N/A |
+-----+

+-----+
| Processes:
|
| GPU   GI    CI        PID   Type   Process name                        GPU
| Memory |
| ID     ID
|
+-----+
| No running processes found
|
+-----+
```

We will use 64 bit floats with JAX in order to increase the precision.

```
jax.config.update("jax_enable_x64", True)
```

## 17.2 Structure

In this section we describe the main features of the model.

### 17.2.1 Output, Consumption and Debt

A small open economy is endowed with an exogenous stochastically fluctuating potential output stream  $\{y_t\}$ .

Potential output is realized only in periods in which the government honors its sovereign debt.

The output good can be traded or consumed.

The sequence  $\{y_t\}$  is described by a Markov process with stochastic density kernel  $p(y, y')$ .

Households within the country are identical and rank stochastic consumption streams according to

$$\mathbb{E} \sum_{t=0}^{\infty} \beta^t u(c_t) \quad (17.1)$$

Here

- $0 < \beta < 1$  is a time discount factor
- $u$  is an increasing and strictly concave utility function

Consumption sequences enjoyed by households are affected by the government's decision to borrow or lend internationally.

The government is benevolent in the sense that its aim is to maximize (17.1).

The government is the only domestic actor with access to foreign credit.

Because households are averse to consumption fluctuations, the government will try to smooth consumption by borrowing from (and lending to) foreign creditors.

### 17.2.2 Asset Markets

The only credit instrument available to the government is a one-period bond traded in international credit markets.

The bond market has the following features

- The bond matures in one period and is not state contingent.
- A purchase of a bond with face value  $B'$  is a claim to  $B'$  units of the consumption good next period.
- To purchase  $B'$  next period costs  $qB'$  now, or, what is equivalent.
- For selling  $-B'$  units of next period goods the seller earns  $-qB'$  of today's goods.
  - If  $B' < 0$ , then  $-qB'$  units of the good are received in the current period, for a promise to repay  $-B'$  units next period.
  - There is an equilibrium price function  $q(B', y)$  that makes  $q$  depend on both  $B'$  and  $y$ .

Earnings on the government portfolio are distributed (or, if negative, taxed) lump sum to households.

When the government is not excluded from financial markets, the one-period national budget constraint is

$$c = y + B - q(B', y)B' \quad (17.2)$$

Here and below, a prime denotes a next period value or a claim maturing next period.

To rule out Ponzi schemes, we also require that  $B \geq -Z$  in every period.

- $Z$  is chosen to be sufficiently large that the constraint never binds in equilibrium.

### 17.2.3 Financial Markets

Foreign creditors

- are risk neutral
- know the domestic output stochastic process  $\{y_t\}$  and observe  $y_t, y_{t-1}, \dots$ , at time  $t$
- can borrow or lend without limit in an international credit market at a constant international interest rate  $r$
- receive full payment if the government chooses to pay
- receive zero if the government defaults on its one-period debt due

When a government is expected to default next period with probability  $\delta$ , the expected value of a promise to pay one unit of consumption next period is  $1 - \delta$ .

Therefore, the discounted expected value of a promise to pay  $B$  next period is

$$q = \frac{1 - \delta}{1 + r} \quad (17.3)$$

Next we turn to how the government in effect chooses the default probability  $\delta$ .

### 17.2.4 Government's Decisions

At each point in time  $t$ , the government chooses between

1. defaulting
2. meeting its current obligations and purchasing or selling an optimal quantity of one-period sovereign debt

Defaulting means declining to repay all of its current obligations.

If the government defaults in the current period, then consumption equals current output.

But a sovereign default has two consequences:

1. Output immediately falls from  $y$  to  $h(y)$ , where  $0 \leq h(y) \leq y$ .
  - It returns to  $y$  only after the country regains access to international credit markets.
1. The country loses access to foreign credit markets.

### 17.2.5 Reentering International Credit Market

While in a state of default, the economy regains access to foreign credit in each subsequent period with probability  $\theta$ .

## 17.3 Equilibrium

Informally, an equilibrium is a sequence of interest rates on its sovereign debt, a stochastic sequence of government default decisions and an implied flow of household consumption such that

1. Consumption and assets satisfy the national budget constraint.
2. The government maximizes household utility taking into account
  - the resource constraint
  - the effect of its choices on the price of bonds
  - consequences of defaulting now for future net output and future borrowing and lending opportunities
1. The interest rate on the government's debt includes a risk-premium sufficient to make foreign creditors expect on average to earn the constant risk-free international interest rate.

To express these ideas more precisely, consider first the choices of the government, which

1. enters a period with initial assets  $B$ , or what is the same thing, initial debt to be repaid now of  $-B$
2. observes current output  $y$ , and
3. chooses either
4. to default, or
5. to pay  $-B$  and set next period's debt due to  $-B'$

In a recursive formulation,

- state variables for the government comprise the pair  $(B, y)$
- $v(B, y)$  is the optimum value of the government's problem when at the beginning of a period it faces the choice of whether to honor or default
- $v_c(B, y)$  is the value of choosing to pay obligations falling due
- $v_d(y)$  is the value of choosing to default

$v_d(y)$  does not depend on  $B$  because, when access to credit is eventually regained, net foreign assets equal 0.

Expressed recursively, the value of defaulting is

$$v_d(y) = u(h(y)) + \beta \int \{\theta v(0, y') + (1 - \theta)v_d(y')\} p(y, y') dy'$$

The value of paying is

$$v_c(B, y) = \max_{B' \geq -Z} \left\{ u(y - q(B', y)B' + B) + \beta \int v(B', y') p(y, y') dy' \right\}$$

The three value functions are linked by

$$v(B, y) = \max\{v_c(B, y), v_d(y)\}$$

The government chooses to default when

$$v_c(B, y) < v_d(y)$$

and hence given  $B'$  the probability of default next period is

$$\delta(B', y) := \int \mathbb{1}\{v_c(B', y') < v_d(y')\} p(y, y') dy' \quad (17.4)$$

Given zero profits for foreign creditors in equilibrium, we can combine (17.3) and (17.4) to pin down the bond price function:

$$q(B', y) = \frac{1 - \delta(B', y)}{1 + r} \quad (17.5)$$

### 17.3.1 Definition of Equilibrium

An *equilibrium* is

- a pricing function  $q(B', y)$ ,
- a triple of value functions  $(v_c(B, y), v_d(y), v(B, y))$ ,
- a decision rule telling the government when to default and when to pay as a function of the state  $(B, y)$ , and
- an asset accumulation rule that, conditional on choosing not to default, maps  $(B, y)$  into  $B'$

such that

- The three Bellman equations for  $(v_c(B, y), v_d(y), v(B, y))$  are satisfied
- Given the price function  $q(B', y)$ , the default decision rule and the asset accumulation decision rule attain the optimal value function  $v(B, y)$ , and
- The price function  $q(B', y)$  satisfies equation (17.5)

## 17.4 Computation

Let's now compute an equilibrium of Arellano's model.

The equilibrium objects are the value function  $v(B, y)$ , the associated default decision rule, and the pricing function  $q(B', y)$ .

We'll use our code to replicate Arellano's results.

After that we'll perform some additional simulations.

We use a slightly modified version of the algorithm recommended by Arellano.

- The appendix to [Arellano, 2008] recommends value function iteration until convergence, updating the price, and then repeating.
- Instead, we update the bond price at every value function iteration step.

The second approach is faster and the two different procedures deliver very similar results.

Here is a more detailed description of our algorithm:

1. Guess a pair of non-default and default value functions  $v_c$  and  $v_d$ .
2. Using these functions, calculate the value function  $v$ , the corresponding default probabilities and the price function  $q$ .
3. At each pair  $(B, y)$ ,
4. update the value of defaulting  $v_d(y)$ .



5. update the value of remaining  $v_c(B, y)$ .
6. Check for convergence. If converged, stop – if not, go to step 2.

We use simple discretization on a grid of asset holdings and income levels.

The output process is discretized using a [quadrature method due to Tauchen](#).

As we have in other places, we accelerate our code using Numba.

We define a namedtuple to store parameters, grids and transition probabilities.

```
ArellanoEconomy = namedtuple('ArellanoEconomy',
    ('β',      # Time discount parameter
     'γ',      # Utility parameter
     'r',      # Lending rate
     'ρ',      # Persistence in the income process
     'η',      # Standard deviation of the income process
     'θ',      # Prob of re-entering financial markets
     'B_size', # Grid size for bonds
     'y_size', # Grid size for income
     'P',      # Markov matrix governing the income process
     'B_grid', # Bond unit grid
     'y_grid', # State values of the income process
     'def_y')) # Default income process
```

```
def create_arellano(B_size=251,      # Grid size for bonds
    B_min=-0.45,      # Smallest B value
    B_max=0.45,      # Largest B value
    y_size=51,      # Grid size for income
    β=0.953,      # Time discount parameter
    γ=2.0,      # Utility parameter
    r=0.017,      # Lending rate
    ρ=0.945,      # Persistence in the income process
    η=0.025,      # Standard deviation of the income process
    θ=0.282,      # Prob of re-entering financial markets
    def_y_param=0.969): # Parameter governing income in default

    # Set up grids
    B_grid = jnp.linspace(B_min, B_max, B_size)
    mc = qe.markov.tauchen(y_size, ρ, η)
    y_grid, P = jnp.exp(mc.state_values), mc.P

    # Put grids on the device
    P = jax.device_put(P)

    # Output received while in default, with same shape as y_grid
    def_y = jnp.minimum(def_y_param * jnp.mean(y_grid), y_grid)

    return ArellanoEconomy(β=β, γ=γ, r=r, ρ=ρ, η=η, θ=θ, B_size=B_size,
        y_size=y_size, P=P,
        B_grid=B_grid, y_grid=y_grid,
        def_y=def_y)
```

Here is the utility function.

```
@jax.jit
def u(c, γ):
    return c** (1-γ) / (1-γ)
```

Here is a function to compute the bond price at each state, given  $v_c$  and  $v_d$ .

```
def compute_q(v_c, v_d, params, sizes, arrays):
    """
    Compute the bond price function  $q(B, y)$  at each  $(B, y)$  pair. The first
    step is to calculate the default probabilities

    
$$\delta(B, y) := \int_{\{y'\} \mid 1\{v_c(B, y') < v_d(y')\}} P(y, y') dy'$$


    """
    # Unpack
    beta, y, r, rho, eta, theta = params
    B_size, y_size = sizes
    P, B_grid, y_grid, def_y = arrays

    # Set up arrays with indices [i_B, i_y, i_yp]
    v_d = jnp.reshape(v_d, (1, 1, y_size))
    v_c = jnp.reshape(v_c, (B_size, 1, y_size))
    P = jnp.reshape(P, (1, y_size, y_size))

    # Compute  $\delta[i_B, i_y]$ 
    default_states = v_c < v_d
    delta = jnp.sum(default_states * P, axis=(2,))

    q = (1 - delta) / (1 + r)
    return q
```

Next we introduce Bellman operators that updated  $v_d$  and  $v_c$ .

```
def T_d(v_c, v_d, params, sizes, arrays):
    """
    The RHS of the Bellman equation when income is at index y_idx and
    the country has chosen to default. Returns an update of v_d.
    """
    # Unpack
    beta, y, r, rho, eta, theta = params
    B_size, y_size = sizes
    P, B_grid, y_grid, def_y = arrays

    B0_idx = jnp.searchsorted(B_grid, 1e-10) # Index at which B is near zero

    current_utility = u(def_y, y)
    v = jnp.maximum(v_c[B0_idx, :], v_d)
    w = theta * v + (1 - theta) * v_d
    A = jnp.reshape(w, (1, y_size))
    cont_value = jnp.sum(A * P, axis=(1,))

    return current_utility + beta * cont_value
```

```
def bellman(v_c, v_d, q, params, sizes, arrays):
    """
    The RHS of the Bellman equation when the country is not in a
    defaulted state on their debt. That is,

    bellman(B, y) =
```

(continues on next page)

(continued from previous page)

```


$$u(y - q(B', y) B' + B) + \beta \mathbb{E}_{\{y'\}} v(B', y') P(y, y')$$


If consumption is not positive then returns -np.inf
"""
# Unpack
β, γ, r, ρ, η, θ = params
B_size, y_size = sizes
P, B_grid, y_grid, def_y = arrays

# Set up c[i_B, i_y, i_Bp]
y_idx = jnp.reshape(jnp.arange(y_size), (1, y_size, 1))
B_idx = jnp.reshape(jnp.arange(B_size), (B_size, 1, 1))
Bp_idx = jnp.reshape(jnp.arange(B_size), (1, 1, B_size))
c = y_grid[y_idx] - q[Bp_idx, y_idx] * B_grid[Bp_idx] + B_grid[B_idx]

# Set up v[i_B, i_y, i_Bp, i_yp] and P[i_B, i_y, i_Bp, i_yp]
v_d = jnp.reshape(v_d, (1, 1, 1, y_size))
v_c = jnp.reshape(v_c, (1, 1, B_size, y_size))
v = jnp.maximum(v_c, v_d)
P = jnp.reshape(P, (1, y_size, 1, y_size))
# Sum over i_yp
continuation_value = jnp.sum(v * P, axis=(3,))

# Return new_v_c[i_B, i_y, i_Bp]
val = jnp.where(c > 0, u(c, γ) + β * continuation_value, -jnp.inf)
return val

```

```

def T_c(v_c, v_d, q, params, sizes, arrays):
    vals = bellman(v_c, v_d, q, params, sizes, arrays)
    return jnp.max(vals, axis=2)

```

```

def get_greedy(v_c, v_d, q, params, sizes, arrays):
    vals = bellman(v_c, v_d, q, params, sizes, arrays)
    return jnp.argmax(vals, axis=2)

```

Let's make JIT-compiled versions of these functions, with the sizes of the arrays declared as static (compile-time constants) in order to help the compiler.

```

compute_q = jax.jit(compute_q, static_argnums=(3,))
T_d = jax.jit(T_d, static_argnums=(3,))
bellman = jax.jit(bellman, static_argnums=(4,))
T_c = jax.jit(T_c, static_argnums=(4,))
get_greedy = jax.jit(get_greedy, static_argnums=(4,))

```

Here is a function that calls these operators in the right sequence.

```

def update_values_and_prices(v_c, v_d, params, sizes, arrays):

    q = compute_q(v_c, v_d, params, sizes, arrays)
    new_v_d = T_d(v_c, v_d, params, sizes, arrays)
    new_v_c = T_c(v_c, v_d, q, params, sizes, arrays)

    return new_v_c, new_v_d

```

We can now write a function that will use an instance of `ArellanoEconomy` and the functions defined above to compute

the solution to our model.

One of the jobs of this function is to take an instance of `ArellanoEconomy`, which is hard for the JIT compiler to handle, and strip it down to more basic objects, which are then passed out to jitted functions.

```
def solve(model, tol=1e-8, max_iter=10_000):
    """
    Given an instance of `ArellanoEconomy`, this function computes the optimal
    policy and value functions.
    """
    # Unpack

     $\beta$ ,  $\gamma$ ,  $r$ ,  $\rho$ ,  $\eta$ ,  $\theta$ , B_size, y_size, P, B_grid, y_grid, def_y = model

    params =  $\beta$ ,  $\gamma$ ,  $r$ ,  $\rho$ ,  $\eta$ ,  $\theta$ 
    sizes = B_size, y_size
    arrays = P, B_grid, y_grid, def_y

     $\beta$ ,  $\gamma$ ,  $r$ ,  $\rho$ ,  $\eta$ ,  $\theta$ , B_size, y_size, P, B_grid, y_grid, def_y = model

    params =  $\beta$ ,  $\gamma$ ,  $r$ ,  $\rho$ ,  $\eta$ ,  $\theta$ 
    sizes = B_size, y_size
    arrays = P, B_grid, y_grid, def_y

    # Initial conditions for v_c and v_d
    v_c = jnp.zeros((B_size, y_size))
    v_d = jnp.zeros((y_size,))

    current_iter = 0
    error = tol + 1
    while (current_iter < max_iter) and (error > tol):
        if current_iter % 100 == 0:
            print(f"Entering iteration {current_iter} with error {error}.")
            new_v_c, new_v_d = update_values_and_prices(v_c, v_d, params,
                                                        sizes, arrays)

            error = jnp.max(jnp.abs(new_v_c - v_c)) + jnp.max(jnp.abs(new_v_d - v_d))
            v_c, v_d = new_v_c, new_v_d
            current_iter += 1

    print(f"Terminating at iteration {current_iter}.")

    q = compute_q(v_c, v_d, params, sizes, arrays)
    B_star = get_greedy(v_c, v_d, q, params, sizes, arrays)
    return v_c, v_d, q, B_star
```

Let's try solving the model.

```
ae = create_arellano()
```

```
2024-08-12 03:21:35.854995: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```

```
%%time
v_c, v_d, q, B_star = solve(ae)
```

```
Entering iteration 0 with error 1.00000001.
```

```
Entering iteration 100 with error 0.017499341639204857.
```

```
Entering iteration 200 with error 0.00014189363558969603.
```

```
Entering iteration 300 with error 1.151467966309383e-06.
```

```
Terminating at iteration 399.
CPU times: user 969 ms, sys: 189 ms, total: 1.16 s
Wall time: 2.37 s
```

We run it again to get rid of compile time.

```
%%time
v_c, v_d, q, B_star = solve(ae)
```

```
Entering iteration 0 with error 1.00000001.
```

```
Entering iteration 100 with error 0.017499341639204857.
```

```
Entering iteration 200 with error 0.00014189363558969603.
```

```
Entering iteration 300 with error 1.151467966309383e-06.
```

```
Terminating at iteration 399.
CPU times: user 396 ms, sys: 153 ms, total: 549 ms
Wall time: 1.42 s
```

Finally, we write a function that will allow us to simulate the economy once we have the policy functions

```
def simulate(model, T, v_c, v_d, q, B_star, key):
    """
    Simulates the Arellano 2008 model of sovereign debt

    Here `model` is an instance of `ArellanoEconomy` and `T` is the length of
    the simulation. Endogenous objects `v_c`, `v_d`, `q` and `B_star` are
    assumed to come from a solution to `model`.

    """
    # Unpack elements of the model
    B_size, y_size = model.B_size, model.y_size
    B_grid, y_grid, P = model.B_grid, model.y_grid, model.P

    B0_idx = jnp.searchsorted(B_grid, 1e-10) # Index at which B is near zero
```

(continues on next page)

(continued from previous page)

```

# Set initial conditions
y_idx = y_size // 2
B_idx = B0_idx
in_default = False

# Create Markov chain and simulate income process
mc = qe.MarkovChain(P, y_grid)
y_sim_indices = mc.simulate_indices(T+1, init=y_idx)

# Allocate memory for outputs
y_sim = jnp.empty(T)
y_a_sim = jnp.empty(T)
B_sim = jnp.empty(T)
q_sim = jnp.empty(T)
d_sim = jnp.empty(T, dtype=int)

# Perform simulation
t = 0
while t < T:

    # Update y_sim and B_sim
    y_sim = y_sim.at[t].set(y_grid[y_idx])
    B_sim = B_sim.at[t].set(B_grid[B_idx])

    # if in default:
    if v_c[B_idx, y_idx] < v_d[y_idx] or in_default:
        # Update y_a_sim
        y_a_sim = y_a_sim.at[t].set(model.def_y[y_idx])
        d_sim = d_sim.at[t].set(1)
        Bp_idx = B0_idx
        # Re-enter financial markets next period with prob 9
        # in_default = False if jnp.random.rand() < model.9 else True
        in_default = False if random.uniform(key) < model.9 else True
        key, _ = random.split(key) # Update the random key
    else:
        # Update y_a_sim
        y_a_sim = y_a_sim.at[t].set(y_sim[t])
        d_sim = d_sim.at[t].set(0)
        Bp_idx = B_star[B_idx, y_idx]

    q_sim = q_sim.at[t].set(q[Bp_idx, y_idx])

    # Update time and indices
    t += 1
    y_idx = y_sim_indices[t]
    B_idx = Bp_idx

return y_sim, y_a_sim, B_sim, q_sim, d_sim

```

## 17.5 Results

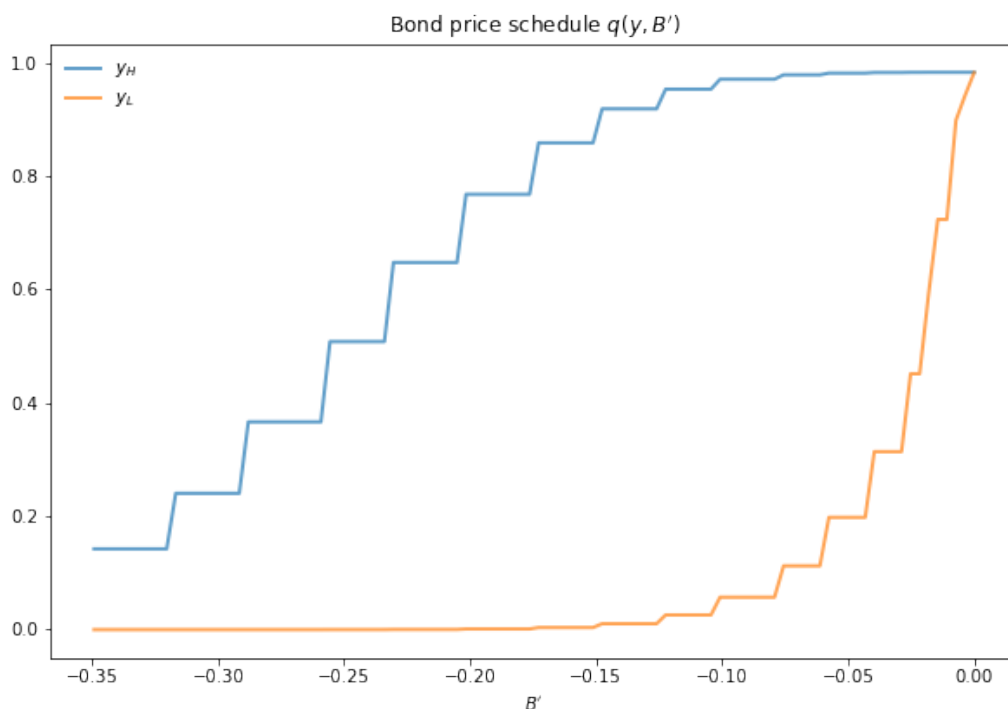
Let's start by trying to replicate the results obtained in [Arellano, 2008].

In what follows, all results are computed using parameter values of `ArellanoEconomy` created by `create_arellano`.

For example,  $r=0.017$  matches the average quarterly rate on a 5 year US treasury over the period 1983–2001.

Details on how to compute the figures are reported as solutions to the exercises.

The first figure shows the bond price schedule and replicates Figure 3 of [Arellano, 2008], where  $y_L$  and  $y_H$  are particular below average and above average values of output  $y$ .



- $y_L$  is 5% below the mean of the  $y$  grid values
- $y_H$  is 5% above the mean of the  $y$  grid values

The grid used to compute this figure was relatively fine (`y_size`, `B_size` = 51, 251), which explains the minor differences between this and Arellano's figure.

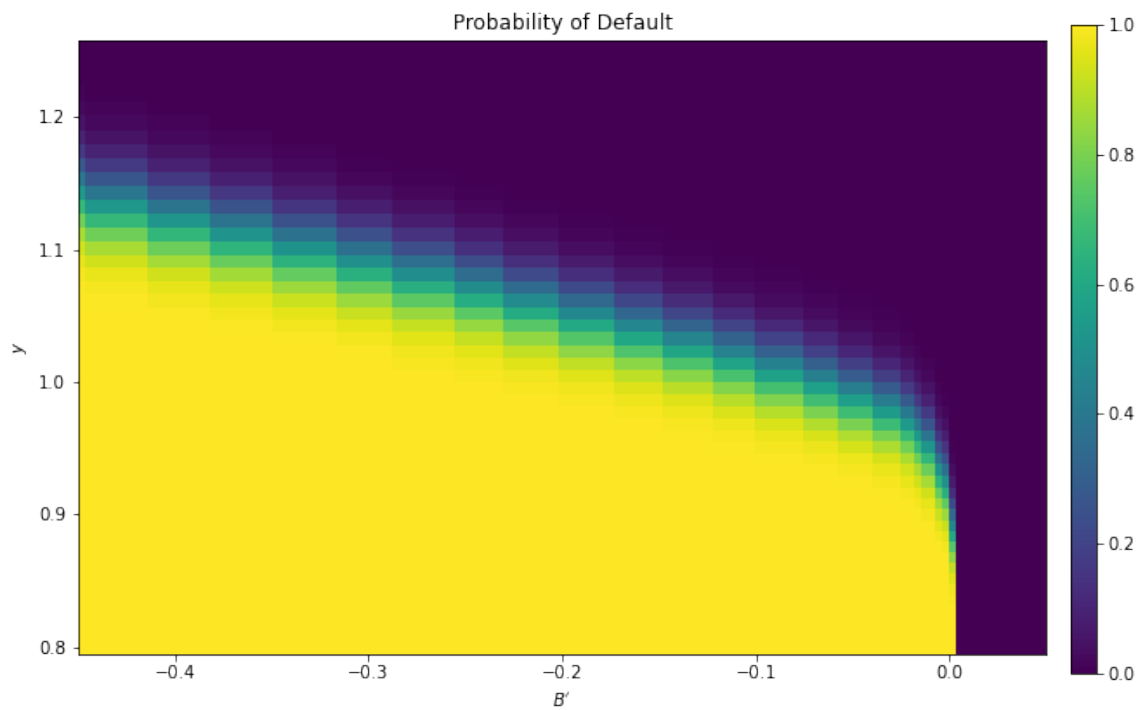
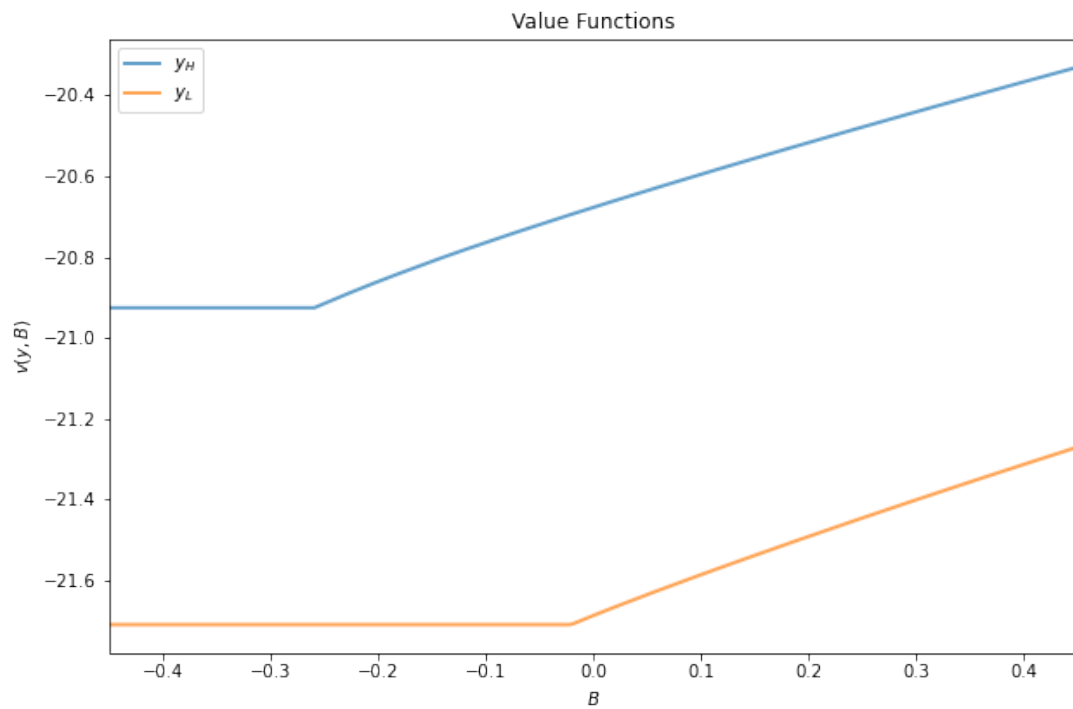
The figure shows that

- Higher levels of debt (larger  $-B'$ ) induce larger discounts on the face value, which correspond to higher interest rates.
- Lower income also causes more discounting, as foreign creditors anticipate greater likelihood of default.

The next figure plots value functions and replicates the right hand panel of Figure 4 of [Arellano, 2008].

We can use the results of the computation to study the default probability  $\delta(B', y)$  defined in (17.4).

The next plot shows these default probabilities over  $(B', y)$  as a heat map.





As anticipated, the probability that the government chooses to default in the following period increases with indebtedness and falls with income.

Next let's run a time series simulation of  $\{y_t\}$ ,  $\{B_t\}$  and  $q(B_{t+1}, y_t)$ .

The grey vertical bars correspond to periods when the economy is excluded from financial markets because of a past default.

One notable feature of the simulated data is the nonlinear response of interest rates.

Periods of relative stability are followed by sharp spikes in the discount rate on government debt.

## 17.6 Exercises

### Exercise 17.6.1

To the extent that you can, replicate the figures shown above

- Use the parameter values listed as defaults in `ArellanoEconomy` created by `create_arellano`.
- The time series will of course vary depending on the shock draws.

### Solution to Exercise 17.6.1

Compute the value function, policy and equilibrium prices

```
ae = create_arellano()
v_c, v_d, q, B_star = solve(ae)
```

```
Entering iteration 0 with error 1.00000001.
```

```
Entering iteration 100 with error 0.017499341639204857.
```

```
Entering iteration 200 with error 0.00014189363558969603.
```

```
Entering iteration 300 with error 1.151467966309383e-06.
```

```
Terminating at iteration 399.
```

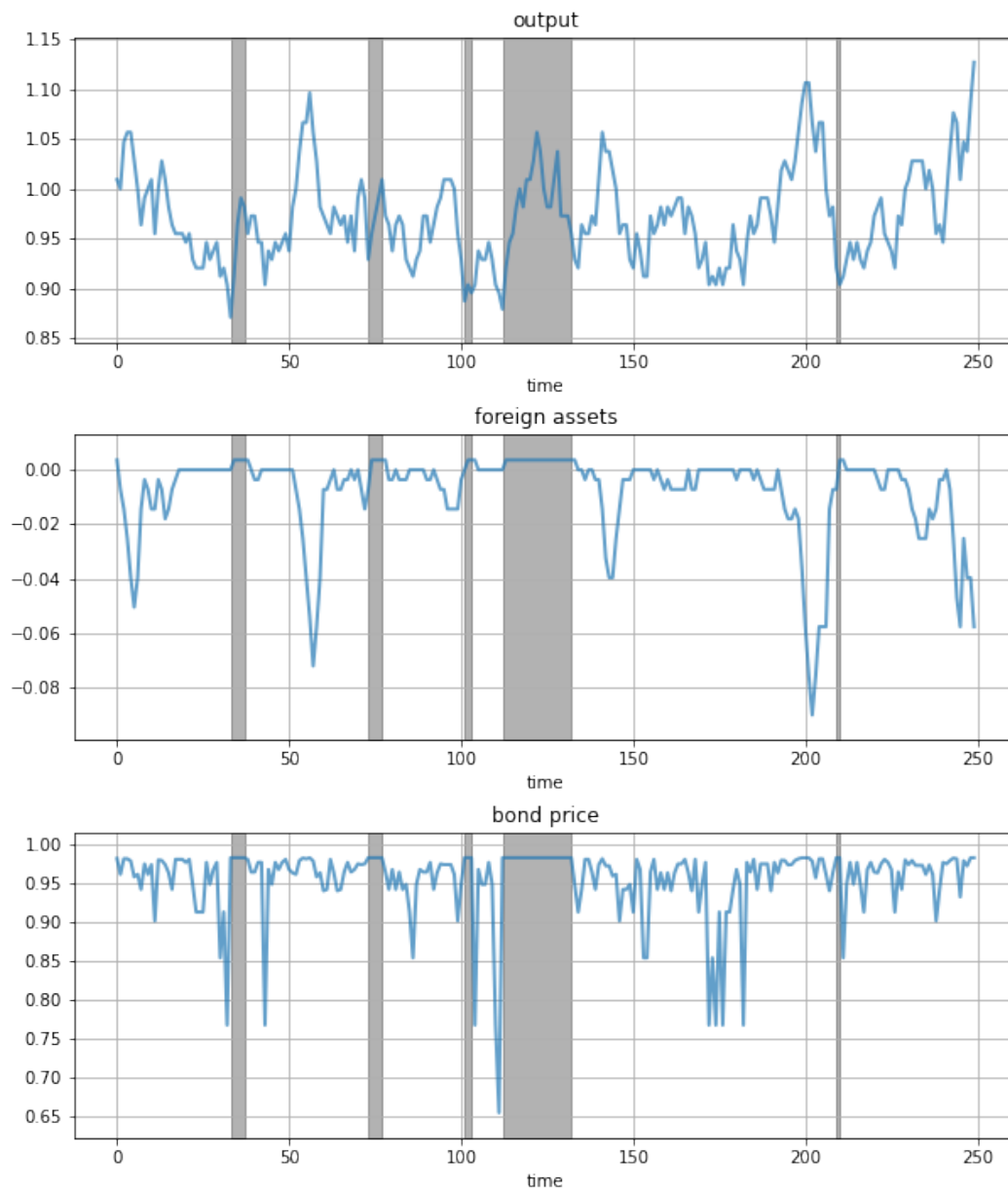
Compute the bond price schedule as seen in figure 3 of [Arellano, 2008]

```
# Unpack some useful names
B_grid, y_grid, P = ae.B_grid, ae.y_grid, ae.P
B_size, y_size = ae.B_size, ae.y_size
r = ae.r

# Create "Y High" and "Y Low" values as 5% devs from mean
high, low = jnp.mean(y_grid) * 1.05, jnp.mean(y_grid) * .95
iy_high, iy_low = (jnp.searchsorted(y_grid, x) for x in (high, low))

fig, ax = plt.subplots(figsize=(10, 6.5))
```

(continues on next page)



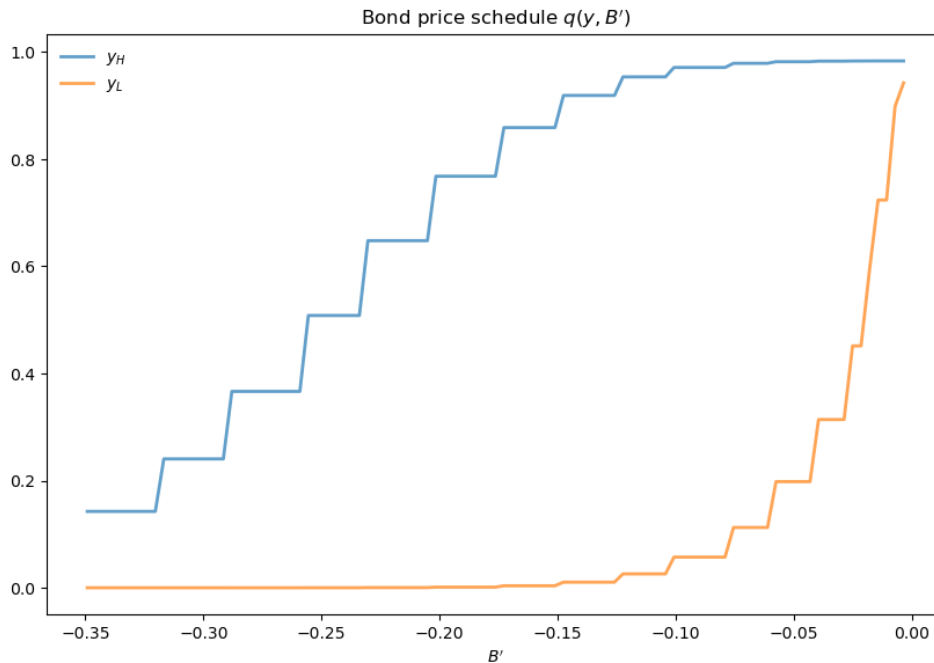
(continued from previous page)

```

ax.set_title("Bond price schedule  $q(y, B')$ ")

# Extract a suitable plot grid
x = []
q_low = []
q_high = []
for i, B in enumerate(B_grid):
    if -0.35 <= B <= 0: # To match fig 3 of Arellano (2008)
        x.append(B)
        q_low.append(q[i, iy_low])
        q_high.append(q[i, iy_high])
ax.plot(x, q_high, label=" $y_H$ ", lw=2, alpha=0.7)
ax.plot(x, q_low, label=" $y_L$ ", lw=2, alpha=0.7)
ax.set_xlabel(" $B'$ ")
ax.legend(loc='upper left', frameon=False)
plt.show()

```



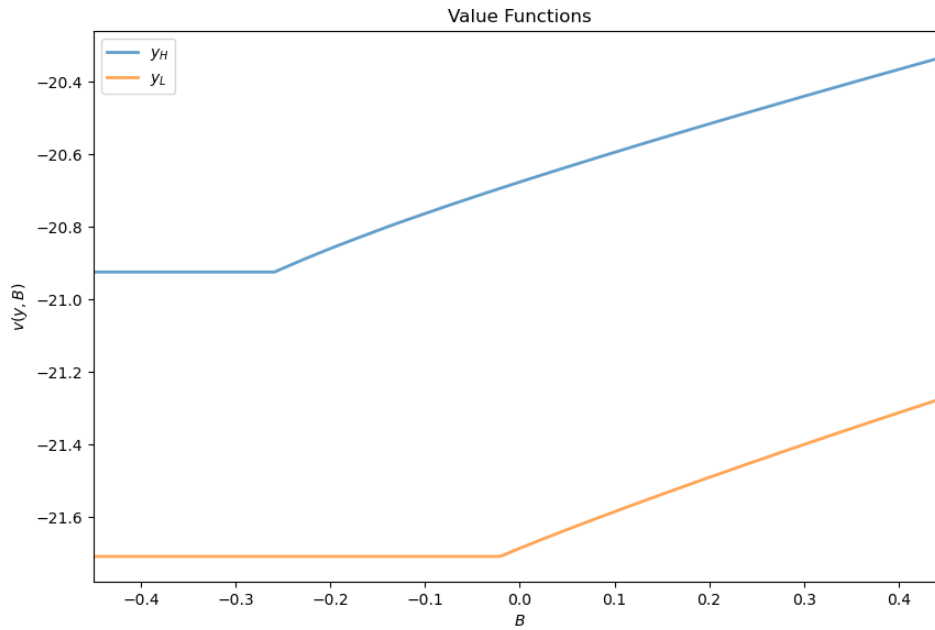
Draw a plot of the value functions

```

v = jnp.maximum(v_c, jnp.reshape(v_d, (1, y_size)))

fig, ax = plt.subplots(figsize=(10, 6.5))
ax.set_title("Value Functions")
ax.plot(B_grid, v[:, iy_high], label=" $y_H$ ", lw=2, alpha=0.7)
ax.plot(B_grid, v[:, iy_low], label=" $y_L$ ", lw=2, alpha=0.7)
ax.legend(loc='upper left')
ax.set_xlabel(" $B$ ", ylabel=" $v(y, B)$ ")
ax.set_xlim(min(B_grid), max(B_grid))
plt.show()

```

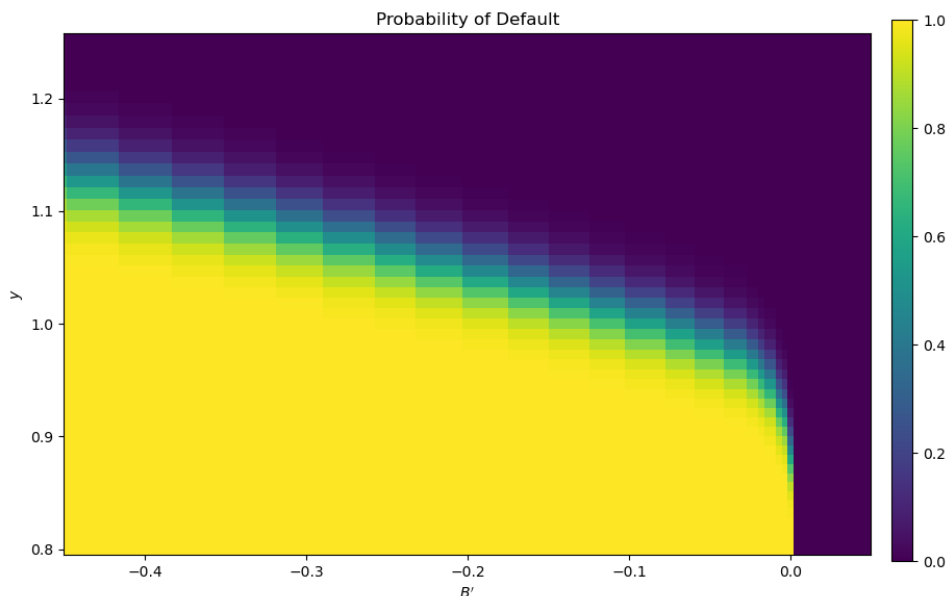


Draw a heat map for default probability

```
# Set up arrays with indices [i_B, i_y, i_yp]
shaped_v_d = jnp.reshape(v_d, (1, 1, y_size))
shaped_v_c = jnp.reshape(v_c, (B_size, 1, y_size))
shaped_P = jnp.reshape(P, (1, y_size, y_size))

# Compute delta[i_B, i_y]
default_states = 1.0 * (shaped_v_c < shaped_v_d)
delta = jnp.sum(default_states * shaped_P, axis=(2,))

# Create figure
fig, ax = plt.subplots(figsize=(10, 6.5))
hm = ax.pcolormesh(B_grid, y_grid, delta.T)
cax = fig.add_axes([.92, .1, .02, .8])
fig.colorbar(hm, cax=cax)
ax.axis([B_grid.min(), 0.05, y_grid.min(), y_grid.max()])
ax.set(xlabel="$B$", ylabel="$y$", title="Probability of Default")
plt.show()
```



Plot a time series of major variables simulated from the model

```
import jax.random as random
T = 250
key = random.PRNGKey(42)
y_sim, y_a_sim, B_sim, q_sim, d_sim = simulate(ae, T, v_c, v_d, q, B_star, key)

# T = 250
# jnp.random.seed(42)
# y_sim, y_a_sim, B_sim, q_sim, d_sim = simulate(ae, T, v_c, v_d, q, B_star)
```

```
# Pick up default start and end dates
start_end_pairs = []
i = 0
while i < len(d_sim):
    if d_sim[i] == 0:
        i += 1
    else:
        # If we get to here we're in default
        start_default = i
        while i < len(d_sim) and d_sim[i] == 1:
            i += 1
        end_default = i - 1
        start_end_pairs.append((start_default, end_default))

plot_series = (y_sim, B_sim, q_sim)
titles = 'output', 'foreign assets', 'bond price'

fig, axes = plt.subplots(len(plot_series), 1, figsize=(10, 12))
fig.subplots_adjust(hspace=0.3)

for ax, series, title in zip(axes, plot_series, titles):
    # Determine suitable y limits
    s_max, s_min = max(series), min(series)
    s_range = s_max - s_min
```

(continues on next page)

(continued from previous page)

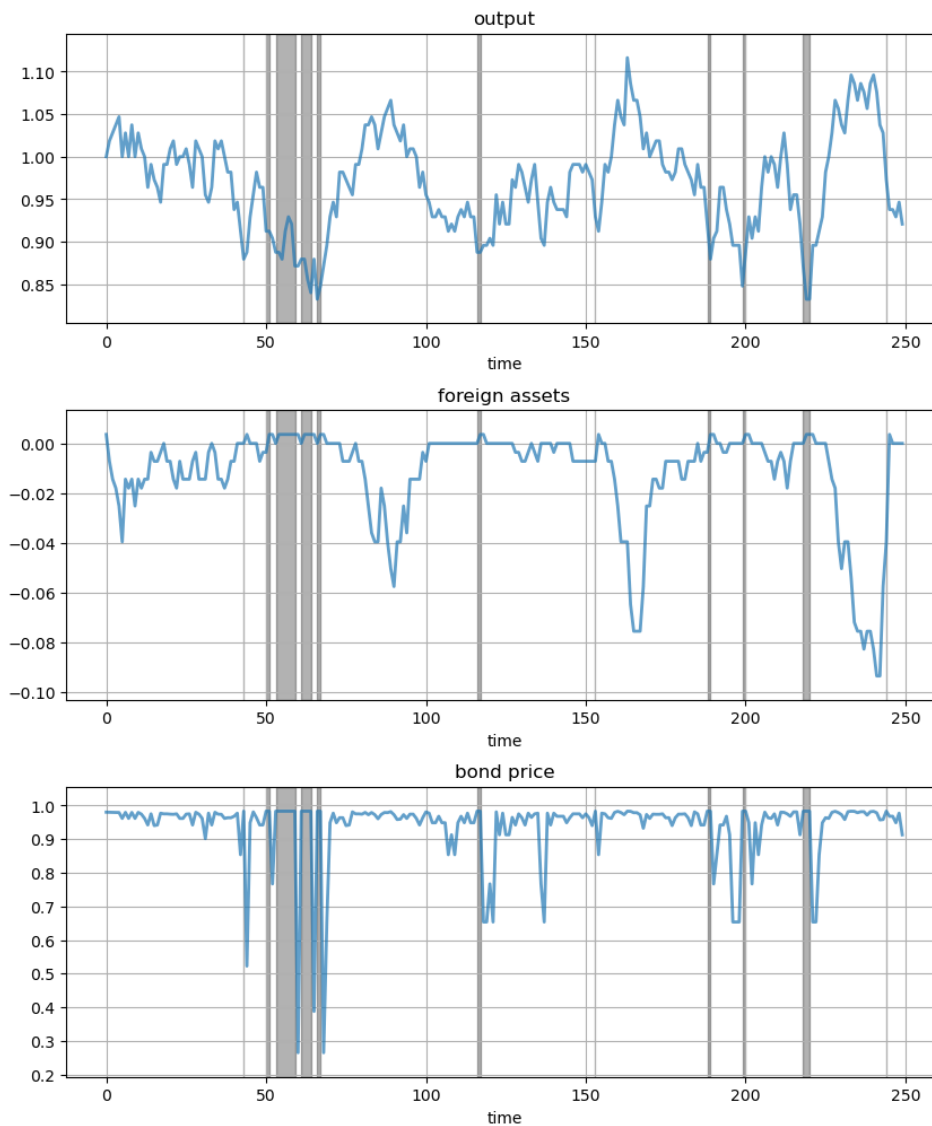
```

y_max = s_max + s_range * 0.1
y_min = s_min - s_range * 0.1
ax.set_ylim(y_min, y_max)
for pair in start_end_pairs:
    ax.fill_between(pair, (y_min, y_min), (y_max, y_max),
                    color='k', alpha=0.3)

ax.grid()
ax.plot(range(T), series, lw=2, alpha=0.7)
ax.set(title=title, xlabel="time")

plt.show()

```



## THE AIYAGARI MODEL

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

## 18.1 Overview

In this lecture, we describe the structure of a class of models that build on work by Truman Bewley [[Bew77](#)].

We begin by discussing an example of a Bewley model due to Rao Aiyagari [[Aiy94](#)].

The model features

- Heterogeneous agents
- A single exogenous vehicle for borrowing and lending
- Limits on amounts individual agents may borrow

The Aiyagari model has been used to investigate many topics, including

- precautionary savings and the effect of liquidity constraints [[Aiy94](#)]
- risk sharing and asset pricing [[HL96](#)]
- the shape of the wealth distribution [[BBZ15](#)]

### 18.1.1 References

The primary reference for this lecture is [[Aiy94](#)].

A textbook treatment is available in chapter 18 of [[LS18](#)].

A less sophisticated version of this lecture (without JAX) can be found [here](#).

## 18.1.2 Preliminaries

We use the following imports

```
import time
import matplotlib.pyplot as plt
import numpy as np
import jax
import jax.numpy as jnp
from collections import namedtuple
```

Let's check the GPU we are running

```
!nvidia-smi
```

```
Mon Aug 12 03:20:23 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5 |
+-----+
| GPU   Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+=====+
|  0  Tesla T4                               On | 00000001:00:00.0 Off |
| Off |
| N/A   30C    P8              9W /  70W |      2MiB / 16384MiB |      0%    |
| Default |
|
| N/A |
+-----+

+-----+
| Processes:
|
| GPU   GI    CI          PID    Type    Process name                        GPU
| Memory |
| ID     ID
|
+=====+
| No running processes found
+-----+
```

We will use 64 bit floats with JAX in order to increase the precision.

```
jax.config.update("jax_enable_x64", True)
```



We will use the following function to compute stationary distributions of stochastic matrices. (For a reference to the algorithm, see p. 88 of [Economic Dynamics](#).)

```
@jax.jit
def compute_stationary(P):
    n = P.shape[0]
    I = jnp.identity(n)
    O = jnp.ones((n, n))
    A = I - jnp.transpose(P) + O
    return jnp.linalg.solve(A, jnp.ones(n))
```

## 18.2 Firms

Firms produce output by hiring capital and labor.

Firms act competitively and face constant returns to scale.

Since returns to scale are constant the number of firms does not matter.

Hence we can consider a single (but nonetheless competitive) representative firm.

The firm's output is

$$Y = AK^\alpha N^{1-\alpha}$$

where

- $A$  and  $\alpha$  are parameters with  $A > 0$  and  $\alpha \in (0, 1)$
- $K$  is aggregate capital
- $N$  is total labor supply (which is constant in this simple version of the model)

The firm's problem is

$$\max_{K, N} \{AK^\alpha N^{1-\alpha} - (r + \delta)K - wN\}$$

The parameter  $\delta$  is the depreciation rate.

These parameters are stored in the following namedtuple.

```
Firm = namedtuple('Firm', ('A', 'N', 'alpha', 'delta'))

def create_firm(A=1.0,
                N=1.0,
                alpha=0.33,
                delta=0.05):
    """
    Create a namedtuple that stores firm data.

    """
    return Firm(A=A, N=N, alpha=alpha, delta=delta)
```

From the first-order condition with respect to capital, the firm's inverse demand for capital is

$$r = A\alpha \left(\frac{N}{K}\right)^{1-\alpha} - \delta \quad (18.1)$$

```
def r_given_k(K, firm):
    """
    Inverse demand curve for capital. The interest rate associated with a
    given demand for capital K.
    """
    A, N, α, δ = firm
    return A * α * (N / K)**(1 - α) - δ
```

Using (18.1) and the firm's first-order condition for labor,

$$w(r) = A(1 - \alpha)(A\alpha/(r + \delta))^{\alpha/(1-\alpha)} \quad (18.2)$$

```
def r_to_w(r, firm):
    """
    Equilibrium wages associated with a given interest rate r.
    """
    A, N, α, δ = firm
    return A * (1 - α) * (A * α / (r + δ))**(α / (1 - α))
```

## 18.3 Households

Infinitely lived households / consumers face idiosyncratic income shocks.

A unit interval of *ex-ante* identical households face a common borrowing constraint.

The savings problem faced by a typical household is

$$\max \mathbb{E} \sum_{t=0}^{\infty} \beta^t u(c_t)$$

subject to

$$a_{t+1} + c_t \leq wz_t + (1 + r)a_t \quad c_t \geq 0, \quad \text{and} \quad a_t \geq -B$$

where

- $c_t$  is current consumption
- $a_t$  is assets
- $z_t$  is an exogenous component of labor income capturing stochastic unemployment risk, etc.
- $w$  is a wage rate
- $r$  is a net interest rate
- $B$  is the maximum amount that the agent is allowed to borrow

The exogenous process  $\{z_t\}$  follows a finite state Markov chain with given stochastic matrix  $P$ .

In this simple version of the model, households supply labor inelastically because they do not value leisure.

Below we provide code to solve the household problem, taking  $r$  and  $w$  as fixed.

### 18.3.1 Primitives and Operators

We will solve the household problem using Howard policy iteration (see Ch 5 of [Dynamic Programming](#)).

First we set up a namedtuple to store the parameters that define a household asset accumulation problem, as well as the grids used to solve it.

```
Household = namedtuple('Household',
                       ('β', 'a_grid', 'z_grid', 'Π'))

def create_household(β=0.96,                # Discount factor
                    Π=[[0.9, 0.1], [0.1, 0.9]], # Markov chain
                    z_grid=[0.1, 1.0],      # Exogenous states
                    a_min=1e-10, a_max=20,   # Asset grid
                    a_size=200):
    """
    Create a namedtuple that stores all data needed to solve the household
    problem, given prices.

    """
    a_grid = jnp.linspace(a_min, a_max, a_size)
    z_grid, Π = map(jnp.array, (z_grid, Π))
    return Household(β=β, a_grid=a_grid, z_grid=z_grid, Π=Π)
```

For now we assume that  $u(c) = \log(c)$ .

(CRRA utility is treated in the exercises.)

```
u = jnp.log
```

Here's a tuple that stores the wage rate and interest rate, as well as a function that creates a price namedtuple with default values.

```
Prices = namedtuple('Prices', ('r', 'w'))

def create_prices(r=0.01, # Interest rate
                 w=1.0): # Wages
    return Prices(r=r, w=w)
```

Now we set up a vectorized version of the right-hand side of the Bellman equation (before maximization), which is a 3D array representing

$$B(a, z, a') = u(wz + (1 + r)a - a') + \beta \sum_{z'} v(a', z') \Pi(z, z')$$

for all  $(a, z, a')$ .

```
@jax.jit
def B(v, household, prices):
    # Unpack
    β, a_grid, z_grid, Π = household
    a_size, z_size = len(a_grid), len(z_grid)
    r, w = prices

    # Compute current consumption as array c[i, j, ip]
    a = jnp.reshape(a_grid, (a_size, 1, 1)) # a[i] -> a[i, j, ip]
```

(continues on next page)

(continued from previous page)

```

z = jnp.reshape(z_grid, (1, z_size, 1))    # z[j]    -> z[i, j, ip]
ap = jnp.reshape(a_grid, (1, 1, a_size))    # ap[ip]  -> ap[i, j, ip]
c = w * z + (1 + r) * a - ap

# Calculate continuation rewards at all combinations of (a, z, ap)
v = jnp.reshape(v, (1, 1, a_size, z_size)) # v[ip, jp] -> v[i, j, ip, jp]
Π = jnp.reshape(Π, (1, z_size, 1, z_size))  # Π[j, jp] -> Π[i, j, ip, jp]
EV = jnp.sum(v * Π, axis=-1)                # sum over last index jp

# Compute the right-hand side of the Bellman equation
return jnp.where(c > 0, u(c) + β * EV, -jnp.inf)

```

The next function computes greedy policies.

```

@jax.jit
def get_greedy(v, household, prices):
    """
    Computes a v-greedy policy σ, returned as a set of indices. If
    σ[i, j] equals ip, then a_grid[ip] is the maximizer at i, j.

    """
    return jnp.argmax(B(v, household, prices), axis=-1) # argmax over ap

```

The following function computes the array  $r_\sigma$  which gives current rewards given policy  $\sigma$ .

```

@jax.jit
def compute_r_σ(σ, household, prices):
    """
    Compute current rewards at each i, j under policy σ. In particular,

        r_σ[i, j] = u((1 + r)a[i] + wz[j] - a'[ip])

    when ip = σ[i, j].

    """
    # Unpack
    β, a_grid, z_grid, Π = household
    a_size, z_size = len(a_grid), len(z_grid)
    r, w = prices

    # Compute r_σ[i, j]
    a = jnp.reshape(a_grid, (a_size, 1))
    z = jnp.reshape(z_grid, (1, z_size))
    ap = a_grid[σ]
    c = (1 + r) * a + w * z - ap
    r_σ = u(c)

    return r_σ

```

The value  $v_\sigma$  of a policy  $\sigma$  is defined as

$$v_\sigma = (I - \beta P_\sigma)^{-1} r_\sigma$$

(See Ch 5 of [Dynamic Programming](#) for notation and background on Howard policy iteration.)

To compute this vector, we set up the linear map  $v \rightarrow R_\sigma v$ , where  $R_\sigma := I - \beta P_\sigma$ .

This map can be expressed as

$$(R_{\sigma}v)(a, z) = v(a, z) - \beta \sum_{z'} v(\sigma(a, z), z') \Pi(z, z')$$

(Notice that  $R_{\sigma}$  is expressed as a linear operator rather than a matrix – this is much easier and cleaner to code, and also exploits sparsity.)

```
@jax.jit
def R_σ(v, σ, household):
    # Unpack
    β, a_grid, z_grid, Π = household
    a_size, z_size = len(a_grid), len(z_grid)

    # Set up the array v[σ[i, j], jp]
    zp_idx = jnp.arange(z_size)
    zp_idx = jnp.reshape(zp_idx, (1, 1, z_size))
    σ = jnp.reshape(σ, (a_size, z_size, 1))
    V = v[σ, zp_idx]

    # Expand Π[j, jp] to Π[i, j, jp]
    Π = jnp.reshape(Π, (1, z_size, z_size))

    # Compute and return v[i, j] - β Σ_jp v[σ[i, j], jp] * Π[j, jp]
    return v - β * jnp.sum(V * Π, axis=-1)
```

The next function computes the lifetime value of a given policy.

```
@jax.jit
def get_value(σ, household, prices):
    """
    Get the lifetime value of policy σ by computing

        v_σ = R_σ^{-1} r_σ

    """
    r_σ = compute_r_σ(σ, household, prices)
    # Reduce R_σ to a function in v
    _R_σ = lambda v: R_σ(v, σ, household)
    # Compute v_σ = R_σ^{-1} r_σ using an iterative routine.
    return jax.scipy.sparse.linalg.bicgstab(_R_σ, r_σ)[0]
```

Here's the Howard policy iteration.

```
def howard_policy_iteration(household, prices,
                           tol=1e-4, max_iter=10_000, verbose=False):
    """
    Howard policy iteration routine.

    """
    β, a_grid, z_grid, Π = household
    a_size, z_size = len(a_grid), len(z_grid)
    σ = jnp.zeros((a_size, z_size), dtype=int)

    v_σ = get_value(σ, household, prices)
    i = 0
    error = tol + 1
```

(continues on next page)

(continued from previous page)

```

while error > tol and i < max_iter:
     $\sigma_{\text{new}}$  = get_greedy( $v_{\sigma}$ , household, prices)
     $v_{\sigma_{\text{new}}}$  = get_value( $\sigma_{\text{new}}$ , household, prices)
    error = jnp.max(jnp.abs( $v_{\sigma_{\text{new}}} - v_{\sigma}$ ))
     $\sigma$  =  $\sigma_{\text{new}}$ 
     $v_{\sigma}$  =  $v_{\sigma_{\text{new}}}$ 
    i = i + 1
    if verbose:
        print(f"Concluded loop {i} with error {error}.")
return  $\sigma$ 

```

As a first example of what we can do, let's compute and plot an optimal accumulation policy at fixed prices.

```

# Create an instance of Household
household = create_household()
prices = create_prices()

```

```

2024-08-12 03:20:24.066778: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.

```

```
r, w = prices
```

```
r, w
```

```
(0.01, 1.0)
```

```
%time  $\sigma_{\text{star}}$  = howard_policy_iteration(household, prices, verbose=True)
```

```

Concluded loop 1 with error 11.366804702491983.
Concluded loop 2 with error 9.574671127216241.
Concluded loop 3 with error 3.9653475941216603.
Concluded loop 4 with error 1.1207182943069665.
Concluded loop 5 with error 0.25239593174329933.
Concluded loop 6 with error 0.12172295828855262.
Concluded loop 7 with error 0.043396115585705886.
Concluded loop 8 with error 0.012131906403222814.
Concluded loop 9 with error 0.005823498169206687.
Concluded loop 10 with error 0.0028609954764462486.
Concluded loop 11 with error 0.0016664747156927007.
Concluded loop 12 with error 0.00041604794390837085.
Concluded loop 13 with error 0.0.
CPU times: user 476 ms, sys: 79.9 ms, total: 555 ms
Wall time: 717 ms

```

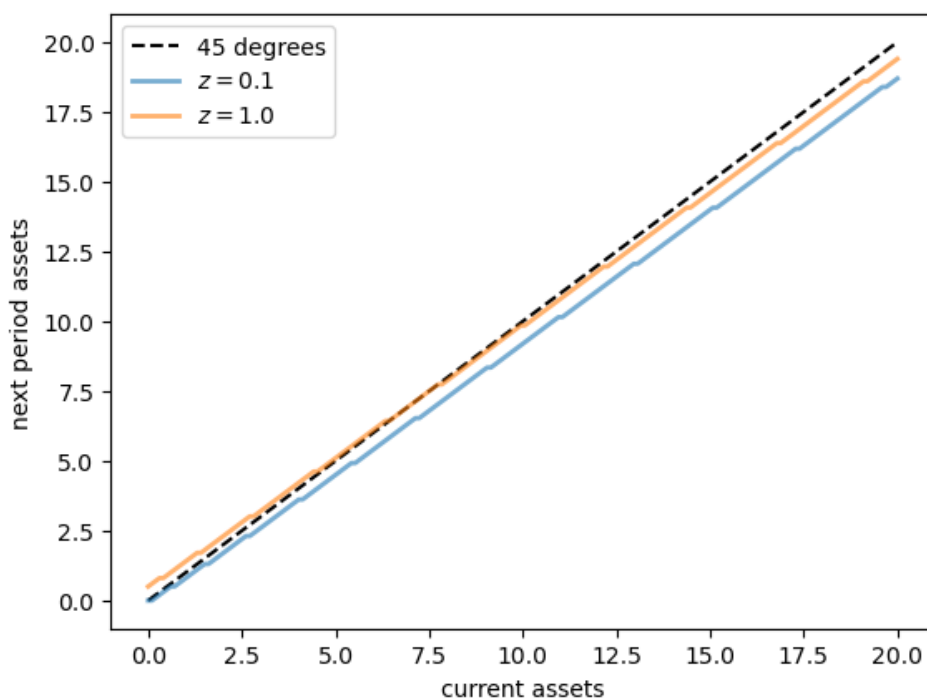
The next plot shows asset accumulation policies at different values of the exogenous state.

```

β, a_grid, z_grid, Π = household

fig, ax = plt.subplots()
ax.plot(a_grid, a_grid, 'k--', label="45 degrees")
for j, z in enumerate(z_grid):
    lb = f'$z = {z:.2}$'
    policy_vals = a_grid[σ_star[:, j]]
    ax.plot(a_grid, policy_vals, lw=2, alpha=0.6, label=lb)
    ax.set_xlabel('current assets')
    ax.set_ylabel('next period assets')
ax.legend(loc='upper left')
plt.show()

```



### 18.3.2 Capital Supply

To start thinking about equilibrium, we need to know how much capital households supply at a given interest rate  $r$ .

This quantity can be calculated by taking the stationary distribution of assets under the optimal policy and computing the mean.

The next function computes the stationary distribution for a given policy  $\sigma$  via the following steps:

- compute the stationary distribution  $\psi = (\psi(a, z))$  of  $P_\sigma$ , which defines the Markov chain of the state  $(a_t, z_t)$  under policy  $\sigma$ .
- sum out  $z_t$  to get the marginal distribution for  $a_t$ .

```

@jax.jit
def compute_asset_stationary(σ, household):
    # Unpack
    β, a_grid, z_grid, Π = household

```

(continues on next page)

(continued from previous page)

```

a_size, z_size = len(a_grid), len(z_grid)

# Construct  $P_\sigma$  as an array of the form  $P_\sigma[i, j, ip, jp]$ 
ap_idx = jnp.arange(a_size)
ap_idx = jnp.reshape(ap_idx, (1, 1, a_size, 1))
 $\sigma$  = jnp.reshape( $\sigma$ , (a_size, z_size, 1, 1))
A = jnp.where( $\sigma$  == ap_idx, 1, 0)
 $\Pi$  = jnp.reshape( $\Pi$ , (1, z_size, 1, z_size))
 $P_\sigma$  = A *  $\Pi$ 

# Reshape  $P_\sigma$  into a matrix
n = a_size * z_size
 $P_\sigma$  = jnp.reshape( $P_\sigma$ , (n, n))

# Get stationary distribution and reshape back onto  $[i, j]$  grid
 $\psi$  = compute_stationary( $P_\sigma$ )
 $\psi$  = jnp.reshape( $\psi$ , (a_size, z_size))

# Sum along the rows to get the marginal distribution of assets
 $\psi_a$  = jnp.sum( $\psi$ , axis=1)
return  $\psi_a$ 

```

Let's give this a test run.

```

 $\psi_a$  = compute_asset_stationary( $\sigma_{\text{star}}$ , household)

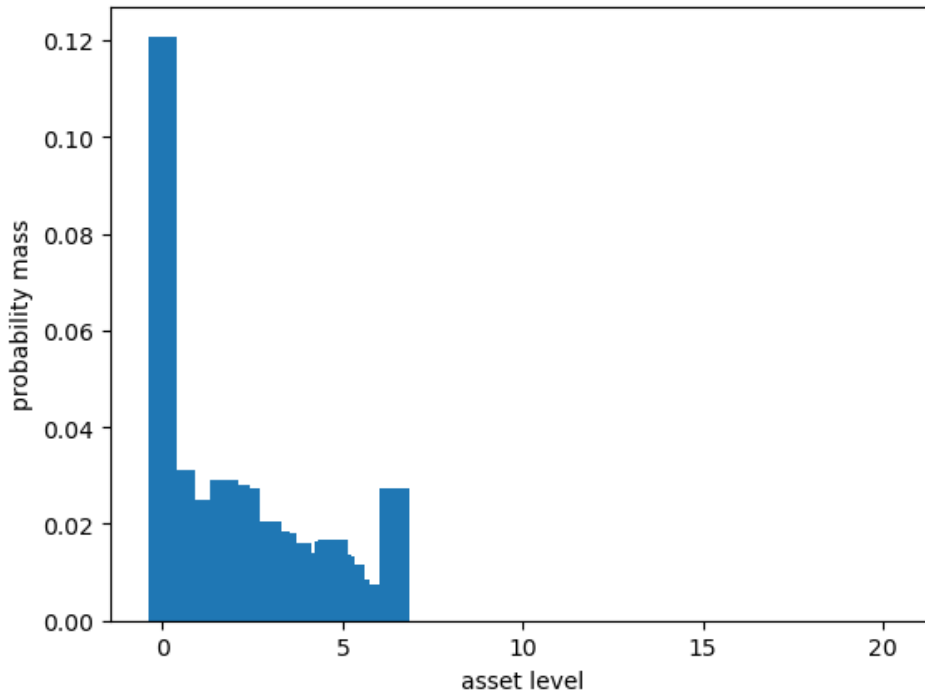
```

```

fig, ax = plt.subplots()
ax.bar(household.a_grid,  $\psi_a$ )
ax.set_xlabel("asset level")
ax.set_ylabel("probability mass")
plt.show()

```





The distribution should sum to one:

```
ψ_a.sum()
```

```
Array(1., dtype=float64)
```

The next function computes aggregate capital supply by households under policy  $\sigma$ , given wages and interest rates.

```
def capital_supply(σ, household):
    """
    Induced level of capital stock under the policy, taking r and w as given.

    """
    β, a_grid, z_grid, Π = household
    ψ_a = compute_asset_stationary(σ, household)
    return float(jnp.sum(ψ_a * a_grid))
```

## 18.4 Equilibrium

We compute a **stationary rational expectations equilibrium** (SREE) as follows:

1. set  $n = 0$ , start with initial guess  $K_0$  for aggregate capital
2. determine prices  $r, w$  from the firm decision problem, given  $K_n$
3. compute the optimal savings policy of the households given these prices
4. compute aggregate capital  $K_{n+1}$  as the mean of steady state capital given this savings policy
5. if  $K_{n+1} \approx K_n$  stop, otherwise go to step 2.

We can write the sequence of operations in steps 2-4 as

$$K_{n+1} = G(K_n)$$

If  $K_{n+1}$  agrees with  $K_n$  then we have a SREE.

In other words, our problem is to find the fixed-point of the one-dimensional map  $G$ .

Here's  $G$  expressed as a Python function:

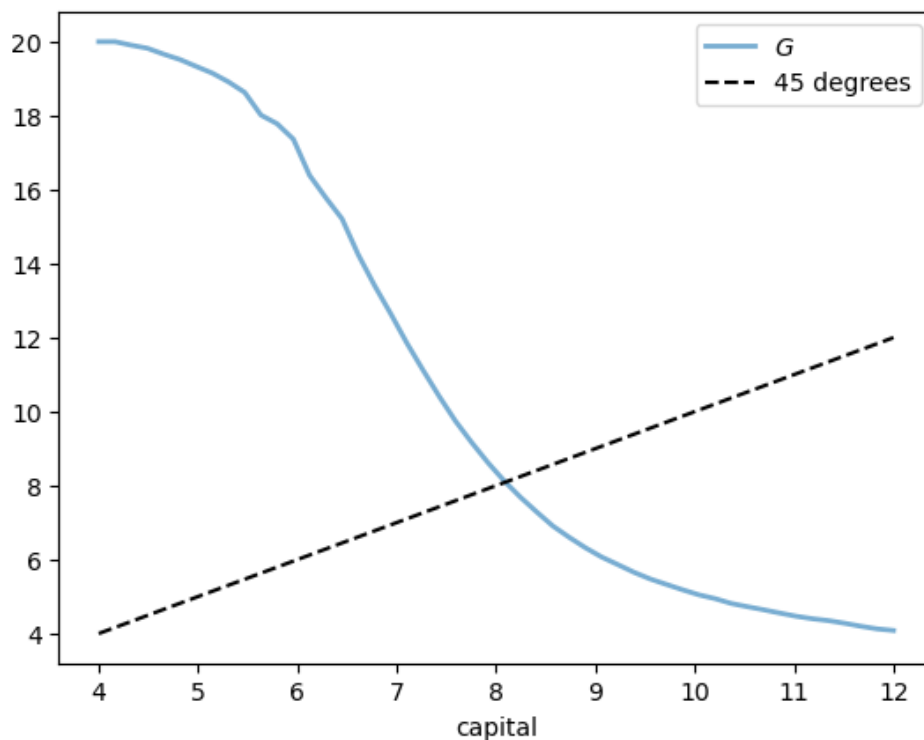
```
def G(K, firm, household):  
    # Get prices r, w associated with K  
    r = r_given_k(K, firm)  
    w = r_to_w(r, firm)  
    # Generate a household object with these prices, compute  
    # aggregate capital.  
    prices = create_prices(r=r, w=w)  
    sigma_star = howard_policy_iteration(household, prices)  
    return capital_supply(sigma_star, household)
```

### 18.4.1 Visual inspection

Let's inspect visually as a first pass.

```
num_points = 50  
firm = create_firm()  
household = create_household()  
k_vals = np.linspace(4, 12, num_points)  
out = [G(k, firm, household) for k in k_vals]
```

```
fig, ax = plt.subplots()  
ax.plot(k_vals, out, lw=2, alpha=0.6, label='$G$')  
ax.plot(k_vals, k_vals, 'k--', label="45 degrees")  
ax.set_xlabel('capital')  
ax.legend()  
plt.show()
```



### 18.4.2 Computing the equilibrium

Now let's compute the equilibrium.

Looking at the figure above, we see that a simple iteration scheme  $K_{n+1} = G(K_n)$  will cycle from high to low values, leading to slow convergence.

As a result, we use a damped iteration scheme of the form

$$K_{n+1} = \alpha K_n + (1 - \alpha)G(K_n)$$

```
def compute_equilibrium(firm, household,
                        K0=6, α=0.99, max_iter=1_000, tol=1e-4,
                        print_skip=10, verbose=False):
    n = 0
    K = K0
    error = tol + 1
    while error > tol and n < max_iter:
        new_K = α * K + (1 - α) * G(K, firm, household)
        error = abs(new_K - K)
        K = new_K
        n += 1
        if verbose and n % print_skip == 0:
            print(f"At iteration {n} with error {error}")
    return K, n
```

```
firm = create_firm()
household = create_household()
print("\nComputing equilibrium capital stock")
```

(continues on next page)

(continued from previous page)

```
start = time.time()
K_star, n = compute_equilibrium(firm, household, K0=6.0, verbose=True)
elapsed = time.time() - start
print(f"Computed equilibrium {K_star:.5} in {n} iterations and {elapsed} seconds")
```

Computing equilibrium capital stock

At iteration 10 with error 0.06490545798825043

At iteration 20 with error 0.03626685465742874

At iteration 30 with error 0.021452703682580676

At iteration 40 with error 0.013383829253038826

At iteration 50 with error 0.008397515264451982

At iteration 60 with error 0.005656370053102933

At iteration 70 with error 0.003750107819943338

At iteration 80 with error 0.0024816210955460605

At iteration 90 with error 0.0017352118525710836

At iteration 100 with error 0.0009559876839571047

At iteration 110 with error 0.0006843412658508186

At iteration 120 with error 0.00046762845631675987

At iteration 130 with error 0.00038193101011607666

At iteration 140 with error 0.00019421626809901227

At iteration 150 with error 0.00017041810881579522

At iteration 160 with error 0.00015412308287032772

At iteration 170 with error 0.00013938615349395889

```
Computed equilibrium 8.0918 in 176 iterations and 40.476239919662476 seconds
```

This is not very fast, given how quickly we can solve the household problem.

You can try varying  $\alpha$ , but usually this parameter is hard to set a priori.

In the exercises below you will be asked to use bisection instead, which generally performs better.

## 18.5 Exercises

### Exercise 18.5.1

Write a new version of `compute_equilibrium` that uses `bisect` from `scipy.optimize` instead of damped iteration.

See if you can make it faster than the previous version.

In `bisect`,

- you should set `xtol=1e-4` to have the same error tolerance as the previous version.
- for the lower and upper bounds of the bisection routine try `a = 1.0` and `b = 20.0`.

### Solution to Exercise 18.5.1

```
from scipy.optimize import bisect
```

We use bisection to find the zero of the function  $h(k) = k - G(k)$ .

```
def compute_equilibrium(firm, household, a=1.0, b=20.0):
    K = bisect(lambda k: k - G(k, firm, household), a, b, xtol=1e-4)
    return K
```

```
firm = create_firm()
household = create_household()
print("\nComputing equilibrium capital stock")
start = time.time()
K_star = compute_equilibrium(firm, household)
elapsed = time.time() - start
print(f"Computed equilibrium capital stock {K_star:.5} in {elapsed} seconds")
```

```
Computing equilibrium capital stock
```

```
Computed equilibrium capital stock 8.0938 in 0.905836820602417 seconds
```

Bisection seems to be faster than the damped iteration scheme.

### Exercise 18.5.2

Show how equilibrium capital stock changes with  $\beta$ .

Use the following values of  $\beta$  and plot the relationship you find.

```
 $\beta$ _vals = np.linspace(0.94, 0.98, 20)
```

---

---

### Solution to Exercise 18.5.2

```
K_vals = np.empty_like( $\beta$ _vals)
K = 6.0 # initial guess

for i,  $\beta$  in enumerate( $\beta$ _vals):
    household = create_household( $\beta$ = $\beta$ )
    K = compute_equilibrium(firm, household, 0.5 * K, 1.5 * K)
    print(f"Computed equilibrium {K:.4} at  $\beta$  = { $\beta$ }")
    K_vals[i] = K
```

```
Computed equilibrium 6.006 at  $\beta$  = 0.94
```

```
Computed equilibrium 6.186 at  $\beta$  = 0.9421052631578947
```

```
Computed equilibrium 6.379 at  $\beta$  = 0.9442105263157894
```

```
Computed equilibrium 6.577 at  $\beta$  = 0.9463157894736841
```

```
Computed equilibrium 6.786 at  $\beta$  = 0.9484210526315789
```

```
Computed equilibrium 7.005 at  $\beta$  = 0.9505263157894737
```

```
Computed equilibrium 7.226 at  $\beta$  = 0.9526315789473684
```

```
Computed equilibrium 7.461 at  $\beta$  = 0.9547368421052631
```

```
Computed equilibrium 7.709 at  $\beta$  = 0.9568421052631578
```

```
Computed equilibrium 7.966 at  $\beta$  = 0.9589473684210525
```

```
Computed equilibrium 8.231 at  $\beta$  = 0.9610526315789474
```

```
Computed equilibrium 8.499 at  $\beta$  = 0.9631578947368421
```

```
Computed equilibrium 8.787 at  $\beta$  = 0.9652631578947368
```

```
Computed equilibrium 9.076 at  $\beta$  = 0.9673684210526315
```

Computed equilibrium 9.378 at  $\beta = 0.9694736842105263$

Computed equilibrium 9.687 at  $\beta = 0.971578947368421$

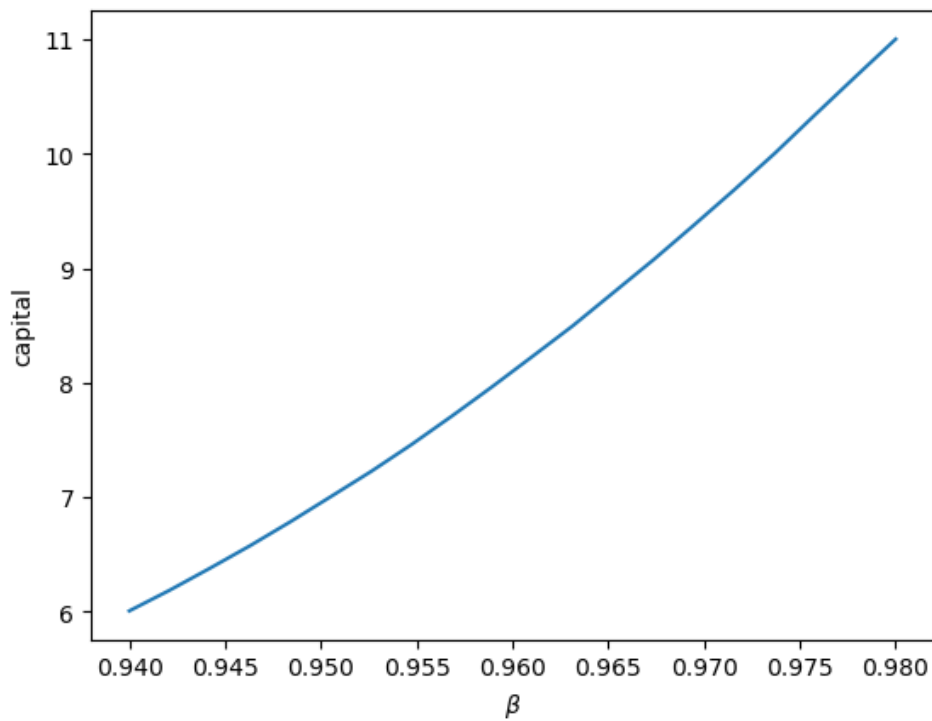
Computed equilibrium 10.0 at  $\beta = 0.9736842105263157$

Computed equilibrium 10.34 at  $\beta = 0.9757894736842105$

Computed equilibrium 10.67 at  $\beta = 0.9778947368421053$

Computed equilibrium 11.0 at  $\beta = 0.98$

```
fig, ax = plt.subplots()
ax.plot( $\beta\_vals$ ,  $K\_vals$ , ms=2)
ax.set_xlabel(r'$\beta$')
ax.set_ylabel('capital')
plt.show()
```







## THE HOPENHAYN ENTRY-EXIT MODEL

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

### 19.1 Outline

The Hopenhayn (1992, ECMA) entry-exit model is a highly influential (partial equilibrium) heterogeneous agent model where

- the agents receiving idiosyncratic shocks are firms, and
- the model produces a non-trivial firm size distribution and a range of predictions for firm dynamics.

We are going to study an extension of the basic model that has unbounded productivity.

Among other things, this allows us to match the heavy tail observed in the firm size distribution data.

In addition to what’s in Anaconda, this lecture will need the following libraries:

```
!pip install quantecon
```

We will use the following imports.

```
import jax.numpy as jnp
import jax
import quantecon as qe
import matplotlib.pyplot as plt
from collections import namedtuple
```

## 19.2 The Model

### 19.2.1 Environment

There is a single good produced by a continuum of competitive firms.

The productivity  $\varphi_t = \varphi_t^i$  of each firm evolves randomly on  $\mathbb{R}_+$ .

- the  $i$  superscript indicates the  $i$ -th firm and we drop it in what follows

A firm with productivity  $\varphi$  facing output price  $p$  and wage rate  $w$  earns current profits

$$\pi(\varphi, p) = p\varphi n^\theta - c - wn$$

Here

- $\theta \in (0, 1)$  is a productivity parameter,
- $c$  is a fixed cost and
- $n$  is labor input.

Maximizing over  $n$  leads to

$$\pi(\varphi, p) = (1 - \theta)(p\varphi)^\eta \left(\frac{\theta}{w}\right)^{\theta/\eta} - c$$

where  $\eta := 1/(1 - \theta)$ .

Output is

$$q(\varphi, p) = \varphi^\eta \left(\frac{p\theta}{w}\right)^{\theta/\eta}$$

Productivity of incumbents is independent across firms and grows according to

$$\varphi_{t+1} \sim \Gamma(\varphi_t, d\varphi')$$

Here  $\Gamma$  is a Markov transition kernel, so that  $\Gamma(\varphi, d\varphi')$  is the distribution over  $\mathbb{R}_+$  given  $\varphi \in \mathbb{R}_+$ .

For example,

$$\int \pi(\varphi', p) \Gamma(\varphi, d\varphi') = \mathbb{E}[\pi(\varphi_{t+1}, p) \mid \varphi_t = \varphi]$$

New entrants obtain their productivity draw from a fixed distribution  $\gamma$ .

Demand is

$$D(p) = 1/p$$

### 19.2.2 Firm decisions

The intertemporal firm decision problem is choosing when to exit.

The timing is:

1. produce and receive current profits  $\pi(\varphi, p)$
2. decide whether to continue or exit.

Scrap value is set to zero and the discount rate is  $\beta \in (0, 1)$ .

The firm makes their stop-continue choice by first solving for the value function, which is the unique solution in a class of continuous functions  $\mathcal{C}$  (details omitted) to

$$v(\varphi, p) = \pi(\varphi, p) + \beta \max \left\{ 0, \int v(\varphi', p) \Gamma(\varphi, d\varphi') \right\}. \quad (19.1)$$

Let  $\bar{v}$  be the unique solution to this functional equation.

Given  $\bar{v}$ , we let  $\bar{\varphi}$  be the exit threshold function defined by

$$\bar{\varphi}(p) := \min \left\{ \varphi \geq 0 \mid \int \bar{v}(\varphi', p) \Gamma(\varphi, d\varphi') \geq 0 \right\}. \quad (19.2)$$

With the convention that incumbents who are indifferent remain rather than exit, an incumbent with productivity  $\varphi$  exits if and only if  $\varphi < \bar{\varphi}(p)$ .

### 19.2.3 Equilibrium

Now we describe a stationary recursive equilibrium in the model, where

- the goods market clears (supply equals demand)
- the firm size distribution is stationary over time (an invariance condition)
- the mass of entrants equals the mass of exiting firms, and
- entering firms satisfy a break-even condition.

The break-even condition means that expected lifetime value given an initial draw from  $\gamma$  is just equal to the fixed cost of entry, which we denote by  $c_e$ .

Let's now state this more formally

To begin, let  $\mathcal{B}$  be the Borel subsets of  $\mathbb{R}_+$  and  $\mathcal{M}$  be all measures on  $\mathcal{B}$ .

Taking  $\bar{v}$  and  $\bar{\varphi}$  as defined in the previous section, a **stationary recursive equilibrium** is a triple

$$(p, M, \mu) \quad \text{in} \quad \mathcal{E} := (0, \infty) \times (0, \infty) \times \mathcal{M},$$

with  $p$  understood as price,  $M$  as mass of entrants, and  $\mu$  as a distribution of firms over productivity levels, such that the goods market clears, or

$$\int q(\varphi, p) \mu(d\varphi) = D(p), \quad (19.3)$$

the *invariance condition* over the firm distribution

$$\mu(B) = \int \Gamma(\varphi, B) \mathbf{1}\{\varphi \geq \bar{\varphi}(p)\} \mu(d\varphi) + M \gamma(B) \text{ for all } B \in \mathcal{B}, \quad (19.4)$$

holds (see below for intuition), the *equilibrium entry condition*

$$\int \bar{v}(\varphi, p) \gamma(d\varphi) = c_e \quad (19.5)$$

holds, and the *balanced entry and exit condition*

$$M = \mu\{\varphi < \bar{\varphi}(p)\} \quad (19.6)$$

is verified.

The invariance condition says that, for any subset  $B$  of the state space, the mass of firms with productivity in set  $B$  today (the left-hand side) is equal to the mass of firms with productivity in set  $B$  tomorrow (the right-hand side).

### 19.2.4 Computing equilibrium

We compute the equilibrium as follows:

For the purposes of this section, we insert balanced entry and exit into the time invariance condition, yielding

$$\mu(B) = \int \Pi(\varphi, B) \mu(d\varphi) \quad \text{for all } B \in \mathcal{B}, \quad (19.7)$$

where  $\Pi$  is the transition kernel on  $\mathbb{R}_+$  defined by

$$\Pi(\varphi, B) = \Gamma(\varphi, B) \mathbf{1}\{\varphi \geq \bar{\varphi}(p)\} + \mathbf{1}\{\varphi < \bar{\varphi}(p)\} \gamma(B). \quad (19.8)$$

Under our assumptions, for each  $p > 0$ , there exists a unique  $\mu$  satisfying this invariance law.

Let  $\mathcal{P}$  be the probability measures on  $\mathbb{R}_+$ .

The unique stationary equilibrium can be computed as follows:

1. Obtain  $\bar{v}$  as the unique solution to the Bellman equation, and then calculate the exit threshold function  $\bar{\varphi}$ .
2. Solve for the equilibrium entry price  $p^*$  by solving the entry condition.
3. Define  $\Pi$  as above, using  $p^*$  as the price, and compute  $\mu$  as the unique solution to the invariance condition.
4. Rescale  $\mu$  by setting  $s := D(p^*) / \int q(\varphi, p^*) \mu(d\varphi)$  and then  $\mu^* := s \mu$ .
5. Obtain the mass of entrants via  $M^* = \mu^* \{\varphi < \bar{\varphi}(p^*)\}$ .

When we compute the distribution  $\mu$  in step 3, we will use the fact that it is the stationary distribution of the Markov transition kernel  $\Pi(\varphi, d\varphi)$ .

This transition kernel turns out to be *ergodic*, which means that if we simulate a cross-section of firms according to  $\Pi$  for a large number of periods, the resulting sample (cross-section of productivities) will approach a set of IID draws from  $\mu$ .

This allows us to

1. compute integrals with respect to the distribution using Monte Carlo, and
2. investigate the shape and properties of the stationary distribution.

### 19.2.5 Specification of dynamics

Before solving the model we need to specify  $\Gamma$  and  $\gamma$ .

We assume  $\Gamma(\varphi, d\varphi')$  is given by

$$\varphi_{t+1} = A_{t+1} \varphi_t$$

We assume that  $(A_t)$  is IID over time, independent across firms, and lognormal  $LN(m_a, \sigma_a)$ .

(This means that incumbents follow Gibrat's law, which is a reasonable assumption for medium to large firms – and hence for incumbents.)

New entrants are drawn from a lognormal distribution  $LN(m_e, \sigma_e)$ .

## 19.3 Code

We use 64 bit floats for extra precision.

```
jax.config.update("jax_enable_x64", True)
```

We store the parameters, grids and Monte Carlo draws in a namedtuple.

```
Parameters = namedtuple("Parameters",
    ("β",          # discount factor
     "θ",          # labor productivity
     "c",          # fixed cost in production
     "c_e",        # entry cost
     "w",          # wages
     "m_a",        # productivity shock location parameter
     "σ_a",        # productivity shock scale parameter
     "m_e",        # new entrant location parameter
     "σ_e"))       # new entrant scale parameter
```

```
Grids = namedtuple("Grids",
    ("ψ_grid",     # productivity grid
     "E_draws",    # entry size draws for Monte Carlo
     "A_draws"))   # productivity shock draws for Monte Carlo
```

```
Model = namedtuple("Model",
    ("parameters", # instance of Parameters
     "grids"))     # instance of Grids
```

```
def create_model(β=0.95,          # discount factor
                 θ=0.3,           # labor productivity
                 c=4.0,           # fixed cost in production
                 c_e=1.0,         # entry cost
                 w=1.0,           # wages
                 m_a=-0.012,      # productivity shock location parameter
                 σ_a=0.1,         # productivity shock scale parameter
                 m_e=1.0,         # new entrant location parameter
                 σ_e=0.2,         # new entrant scale parameter
                 ψ_grid_max=5,    # productivity grid max
                 ψ_grid_size=100, # productivity grid size
                 E_draw_size=200, # entry MC integration size
                 A_draw_size=200, # prod shock MC integration size
                 seed=1234):      # Seed for MC draws
    """
    Create an instance of the `namedtuple` Model using default parameter values.
    """

    # Test stability
    assert m_a + σ_a**2 / (2 * (1 - θ)) < 0, "Stability condition fails"
    # Build grids and initialize random number generator
    ψ_grid = jnp.linspace(0, ψ_grid_max, ψ_grid_size)
    key, subkey = jax.random.split(jax.random.PRNGKey(seed))
    # Generate a sample of draws of A for Monte Carlo integration
    A_draws = jnp.exp(m_a + σ_a * jax.random.normal(key, (A_draw_size,)))
    # Generate a sample of draws from γ for Monte Carlo
```

(continues on next page)

(continued from previous page)

```

E_draws = jnp.exp(m_e +  $\sigma_e$  * jax.random.normal(subkey, (E_draw_size,)))
# Build namedtuple and return
parameters = Parameters( $\beta$ ,  $\theta$ , c, c_e, w, m_a,  $\sigma_a$ , m_e,  $\sigma_e$ )
grids = Grids( $\psi$ _grid, E_draws, A_draws)
model = Model(parameters, grids)
return model

```

Let us write down functions for profits and output.

```

@jax.jit
def  $\pi$ ( $\psi$ , p, parameters):
    """ Profits. """
    # Unpack
     $\beta$ ,  $\theta$ , c, c_e, w, m_a,  $\sigma_a$ , m_e,  $\sigma_e$  = parameters
    # Compute profits
    return (1 -  $\theta$ ) * (p *  $\psi$ )**(1/(1 -  $\theta$ )) * ( $\theta$ /w)**( $\theta$ /(1 -  $\theta$ )) - c

```

```

@jax.jit
def q( $\psi$ , p, parameters):
    """ Output. """
    # Unpack
     $\beta$ ,  $\theta$ , c, c_e, w, m_a,  $\sigma_a$ , m_e,  $\sigma_e$  = parameters
    # Compute output
    return  $\psi$ ** (1/(1 -  $\theta$ )) * (p *  $\theta$ /w)**( $\theta$ /(1 -  $\theta$ ))

```

Let's write code to simulate a cross-section of firms given a particular value for the exit threshold (rather than an exit threshold function).

Firms that exit are immediately replaced by a new entrant, drawn from  $\gamma$ .

Our first function updates by one step

```

def update_cross_section( $\psi$ _bar,  $\psi$ _vec, key, parameters, num_firms):
    # Unpack
     $\beta$ ,  $\theta$ , c, c_e, w, m_a,  $\sigma_a$ , m_e,  $\sigma_e$  = parameters
    # Update
    Z = jax.random.normal(key, (2, num_firms)) # long rows for row-major arrays
    incumbent_draws =  $\psi$ _vec * jnp.exp(m_a +  $\sigma_a$  * Z[0, :])
    new_firm_draws = jnp.exp(m_e +  $\sigma_e$  * Z[1, :])
    return jnp.where( $\psi$ _vec >=  $\psi$ _bar, incumbent_draws, new_firm_draws)

```

```
update_cross_section = jax.jit(update_cross_section, static_argnums=(4,))
```

Our next function runs the cross-section forward in time `sim_length` periods.

```

def simulate_firms( $\psi$ _bar, parameters, grids,
                  sim_length=200, num_firms=1_000_000, seed=12):
    """
    Simulate a cross-section of firms when the exit threshold is  $\psi$ _bar.
    """
    # Set initial conditions to the threshold value
     $\psi$ _vec = jnp.ones((num_firms,)) *  $\psi$ _bar
    key = jax.random.PRNGKey(seed)

```

(continues on next page)

(continued from previous page)

```

# Iterate forward in time
for t in range(sim_length):
    key, subkey = jax.random.split(key)
    psi_vec = update_cross_section(psi_bar, psi_vec, subkey, parameters, num_firms)
return psi_vec

```

Here's a utility function to compute the expected value

$$\int v(\varphi') \Gamma(\varphi, d\varphi') = \mathbb{E}v(A_{t+1}\varphi)$$

given  $\varphi$

```

@jax.jit
def _compute_exp_value_at_phi(v, psi, grids):
    """
    Compute

         $E[v(\psi') | \psi] = E_v(A \psi)$ 

    using linear interpolation and Monte Carlo.
    """
    # Unpack
    psi_grid, E_draws, A_draws = grids
    # Set up V
    Apsi = A_draws * psi
    vApsi = jnp.interp(Apsi, psi_grid, v) # v(A_j psi) for all j
    # Return mean
    return jnp.mean(vApsi) # (1/n) \sum_j v(A_j psi)

```

Now let's vectorize this function in  $\varphi$  and then write another function that computes the expected value across all  $\varphi$  in  $\psi\_grid$

```

compute_exp_value_at_phi = jax.vmap(_compute_exp_value_at_phi, (None, 0, None))

```

```

@jax.jit
def compute_exp_value(v, grids):
    """
    Compute

         $E[v(\psi_{\text{prime}}) | \psi] = E_v(A \psi)$  for all  $\psi$ , as a vector

    """
    # Unpack
    psi_grid, E_draws, A_draws = grids
    return compute_exp_value_at_phi(v, psi_grid, grids)

```

Here is the Bellman operator  $T$ .

```

@jax.jit
def T(v, p, parameters, grids):
    """ Bellman operator. """
    # Unpack
    beta, theta, c, c_e, w, m_a, sigma_a, m_e, sigma_e = parameters
    psi_grid, E_draws, A_draws = grids

```

(continues on next page)

(continued from previous page)

```

# Compute Tv
EvAψ = compute_exp_value(v, grids)
return π(ψ_grid, p, parameters) + β * jnp.maximum(0.0, EvAψ)

```

The next function takes  $v, p$  as inputs and, conditional on the value function  $v$ , computes the value  $\bar{\varphi}(p)$  that corresponds to the exit value.

```

@jax.jit
def get_threshold(v, grids):
    """ Compute the exit threshold. """
    # Unpack
    ψ_grid, E_draws, A_draws = grids
    # Compute exit threshold: ψ such that E v(A ψ) = 0
    EvAψ = compute_exp_value(v, grids)
    i = jnp.searchsorted(EvAψ, 0.0)
    return ψ_grid[i]

```

We use value function iteration (VFI) to compute the value function  $\bar{v}(\cdot, p)$ , taking  $p$  as given.

VFI is relatively cheap and simple in this setting.

```

@jax.jit
def vfi(p, v_init, parameters, grids, tol=1e-6, max_iter=10_000):
    """
    Implement value function iteration to solve for the value function.
    """
    # Unpack
    ψ_grid, E_draws, A_draws = grids
    # Set up
    def cond_function(state):
        i, v, error = state
        return jnp.logical_and(i < max_iter, error > tol)

    def body_function(state):
        i, v, error = state
        new_v = T(v, p, parameters, grids)
        error = jnp.max(jnp.abs(v - new_v))
        i += 1
        return i, new_v, error

    # Loop till convergence
    init_state = 0, v_init, tol + 1
    state = jax.lax.while_loop(cond_function, body_function, init_state)
    i, v, error = state
    return v

```

```

@jax.jit
def compute_net_entry_value(p, v_init, parameters, grids):
    """
    Returns the net value of entry, which is

    \int v_bar(ψ, p) γ(d ψ) - c_e

    This is the break-even condition for new entrants. The argument
    v_init is used as an initial condition when computing v_bar for VFI.
    """

```

(continues on next page)



(continued from previous page)

```

"""
c_e = parameters.c_e
psi_grid = grids.psi_grid
E_draws = grids.E_draws
v_bar = vfi(p, v_init, parameters, grids)
v_psi = jnp.interp(E_draws, psi_grid, v_bar)
Ev_psi = jnp.mean(v_psi)
return Ev_psi - c_e, v_bar

```

We need to solve for the equilibrium price, which is the  $p$  satisfying

$$\int \bar{v}(\varphi', p) \gamma(d\varphi') = c_e$$

At each price  $p$ , we need to recompute  $\bar{v}(\cdot, p)$  and then take the expectation.

The technique we will use is bisection.

We will write our own bisection routine because, when we shift to a new price, we want to update the initial condition for value function iteration to the value function from the previous price.

```

def compute_p_star(parameters, grids, p_min=1.0, p_max=2.0, tol=10e-5):
    """
    Compute the equilibrium entry price p^* via bisection.

    Return both p^* and the corresponding value function v_bar, which is
    computed as a byproduct.

    Implements the bisection root finding algorithm to find p_star

    """
    psi_grid, E_draws, A_draws = grids
    lower, upper = p_min, p_max
    v_bar = jnp.zeros_like(psi_grid) # Initial condition at first price guess

    while upper - lower > tol:
        mid = 0.5 * (upper + lower)
        entry_val, v_bar = compute_net_entry_value(mid, v_bar, parameters, grids)
        if entry_val > 0: # Root is between lower and mid
            lower, upper = lower, mid
        else: # Root is between mid and upper
            lower, upper = mid, upper

    p_star = 0.5 * (upper + lower)
    return p_star, v_bar

```

We are now ready to compute all elements of the stationary recursive equilibrium.

```

def compute_equilibrium_prices_and_quantities(model):
    """
    Compute

    1. The equilibrium outcomes for p^*, v^* and psi^*, where psi^* is the
       equilibrium exit threshold psi_bar(p^*).
    1. The scaling factor necessary to convert the stationary probability
       distribution mu into the equilibrium firm distribution mu^* = s mu.
    """

```

(continues on next page)

(continued from previous page)

```

2. The equilibrium mass of entrants  $M^* = \mu\{\psi < \psi^*\}$ 

"""
# Unpack
parameters, grids = model
# Compute prices and values
p_star, v_bar = compute_p_star(parameters, grids)
# Get  $\psi_{\text{star}} = \psi_{\text{bar}}(p_{\text{star}})$ , the equilibrium exit threshold
psi_star = get_threshold(v_bar, grids)
# Generate an array of draws from  $\mu$ , the normalized stationary distribution.
psi_sample = simulate_firms(psi_star, parameters, grids)
# Compute  $s$  to scale  $\mu$ 
demand = 1 / p_star
pre_normalized_supply = jnp.mean(q(psi_sample, p_star, parameters))
s = demand / pre_normalized_supply
# Compute  $M^* = \mu\{\psi < \psi_{\text{star}}\}$ 
m_star = s * jnp.mean(psi_sample < psi_star)
# return computed objects
return p_star, v_bar, psi_star, psi_sample, s, m_star

```

## 19.4 Solving the model

### 19.4.1 Preliminary calculations

Let's create an instance of the model.

```

model = create_model()
parameters, grids = model

```

```

2024-08-12 03:22:03.440710: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.

```

Let's see how long it takes to compute the value function at a given price from a cold start.

```

p = 2.0
v_init = jnp.zeros_like(grids.psi_grid) # Initial condition
%time v_bar = vfi(p, v_init, parameters, grids).block_until_ready()

```

```

CPU times: user 290 ms, sys: 6.98 ms, total: 297 ms
Wall time: 345 ms

```

Let's run the code again to eliminate compile time.

```

%time v_bar = vfi(p, v_init, parameters, grids).block_until_ready()

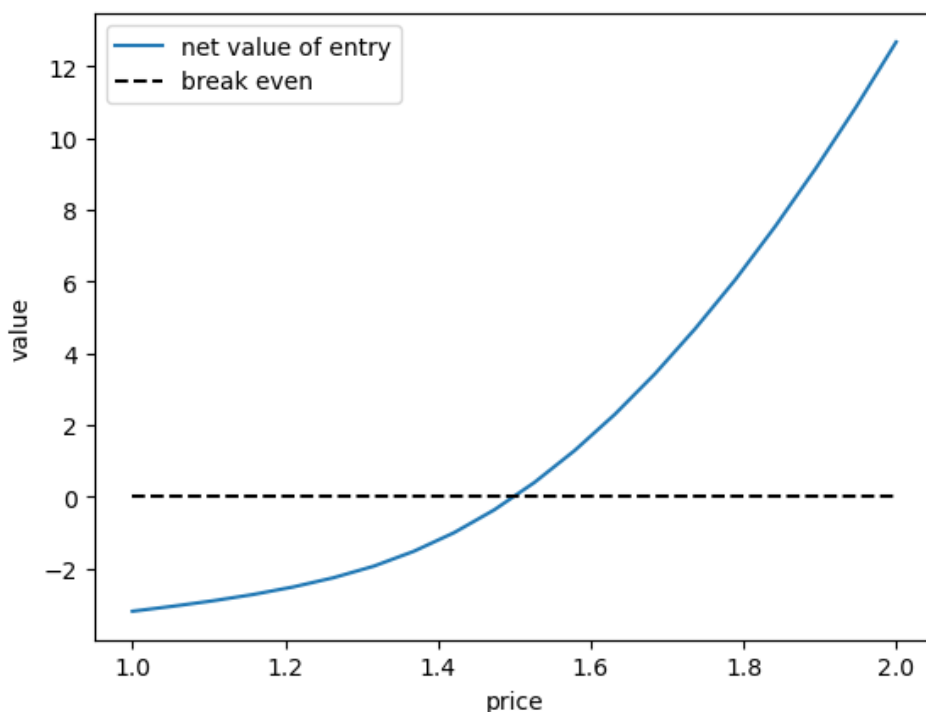
```

```
CPU times: user 28.7 ms, sys: 39 µs, total: 28.7 ms
Wall time: 28.1 ms
```

Let's have a look at the net entry value as a function of price

The root is the equilibrium price at the given parameters

```
p_min, p_max, p_size = 1.0, 2.0, 20
p_vec = jnp.linspace(p_min, p_max, p_size)
entry_vals = []
v_bar = jnp.zeros_like(grids.ψ_grid) # Initial condition at first price guess
for i, p in enumerate(p_vec):
    entry_val, v_bar = compute_net_entry_value(p, v_bar, parameters, grids)
    entry_vals.append(entry_val)
fig, ax = plt.subplots()
ax.plot(p_vec, entry_vals, label="net value of entry")
ax.plot(p_vec, jnp.zeros_like(p_vec), 'k', ls='--', label="break even")
ax.legend()
ax.set_xlabel("price")
ax.set_ylabel("value")
plt.show()
```



Below we solve for the zero of this function to calculate  $p^*$ .

From the figure it looks like  $p^*$  will be close to 1.5.

## 19.4.2 Computing the equilibrium

Now let's try computing the equilibrium

```
%%time
p_star, v_bar, psi_star, psi_sample, s, m_star = \
    compute_equilibrium_prices_and_quantities(model)
```

```
CPU times: user 1.65 s, sys: 96.3 ms, total: 1.75 s
Wall time: 2.17 s
```

Let's run the code again to get rid of compile time.

```
%%time
p_star, v_bar, psi_star, psi_sample, s, m_star = \
    compute_equilibrium_prices_and_quantities(model)
```

```
CPU times: user 281 ms, sys: 10.8 ms, total: 292 ms
Wall time: 671 ms
```

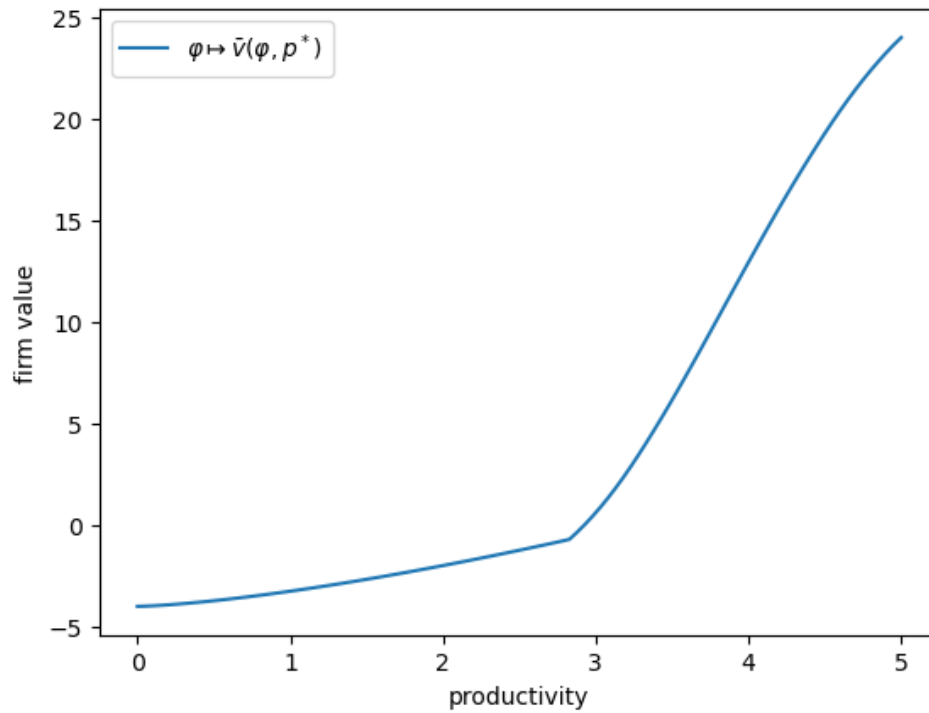
Let's check that  $p^*$  is close to 1.5

```
p_star
```

```
1.500213623046875
```

Here is a plot of the value function  $\bar{v}(\cdot, p^*)$ .

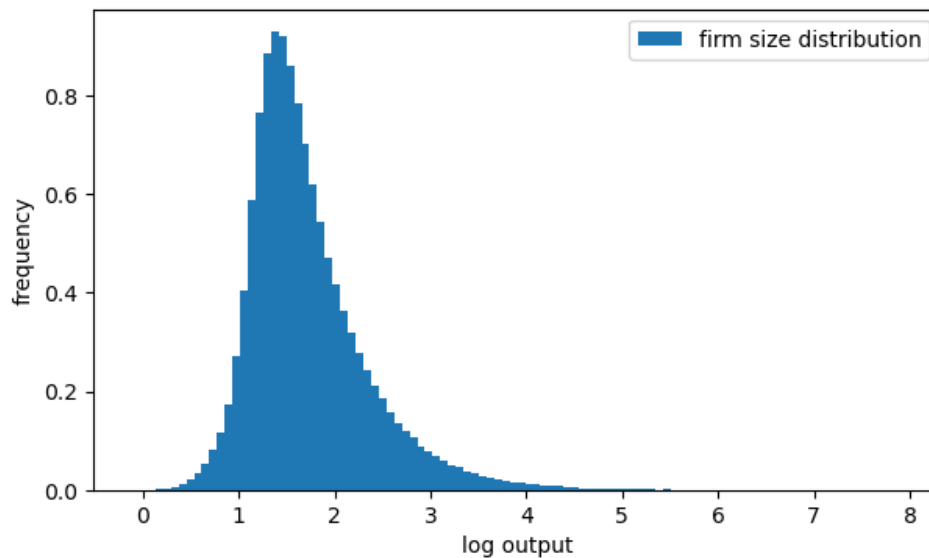
```
fig, ax = plt.subplots()
ax.plot(grids.psi_grid, v_bar, label=r'$\varphi \mapsto \bar{v}(\varphi, p^*)$')
ax.set_xlabel("productivity")
ax.set_ylabel("firm value")
ax.legend()
plt.show()
```



Let's have a look at the firm size distribution, with firm size measured by output.

```
output_dist = q(\(\psi\_sample, p\_star, parameters\))
```

```
fig, ax = plt.subplots(figsize=(7, 4))
ax.hist(jnp.log(output_dist), bins=100, density=True,
        label="firm size distribution")
ax.set_xlabel("log output")
ax.set_ylabel("frequency")
ax.legend()
plt.show()
```



## 19.5 Pareto tails

The firm size distribution shown above appears to have a long right tail.

This matches the observed firm size distribution.

In fact the firm size distribution obeys a **power law**.

More mathematically, the distribution of firm size has a Pareto right hand tail, so that there exist constants  $k, \alpha > 0$  with

$$\mu((x, \infty)) \approx kx^{-\alpha} \text{ when } x \text{ is large}$$

Here  $\alpha$  is called the tail index.

Does the model replicate this feature?

One option is to look at the empirical counter CDF (cumulative distribution).

The idea is as follows: The counter CDF of a random variable  $X$  is

$$G(x) := \mathbb{P}\{X > x\}$$

In the case of a Pareto tailed distribution we have  $\mathbb{P}\{X > x\} \approx kx^{-\alpha}$  for large  $x$ .

Hence, for large  $x$ ,

$$\ln G(x) \approx \ln k - \alpha \ln x$$

The empirical counterpart of  $G$  given sample  $X_1, \dots, X_n$  is

$$G_n(x) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i > x\}$$

For large  $k$  (implying  $G_n \approx G$ ) and large  $x$ , we expect that, for a Pareto-tailed sample,  $\ln G_n$  is approximately linear.

Here's a function to compute the empirical counter CDF:

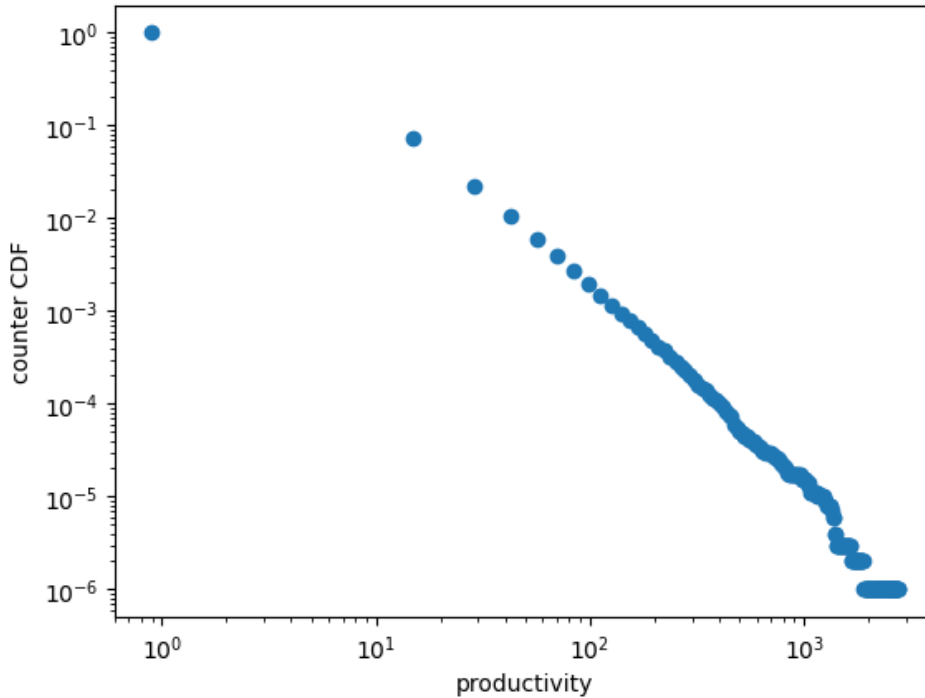
```
def ECCDF(data):
    """
    Return a function that implements the ECCDF given the data.
    """
    def eccdf(x):
        return jnp.mean(data > x)
    return eccdf
```

Let's plot the empirical counter CDF of the output distribution.

```
eccdf = ECCDF(output_dist)

epsilon = 10e-10
x_grid = jnp.linspace(output_dist.min() + epsilon, output_dist.max() - epsilon, 200)
y = [eccdf(x) for x in x_grid]

fig, ax = plt.subplots()
ax.loglog(x_grid, y, 'o', label="ECCDF")
ax.set_xlabel("productivity")
ax.set_ylabel("counter CDF")
plt.show()
```



## 19.6 Exercise

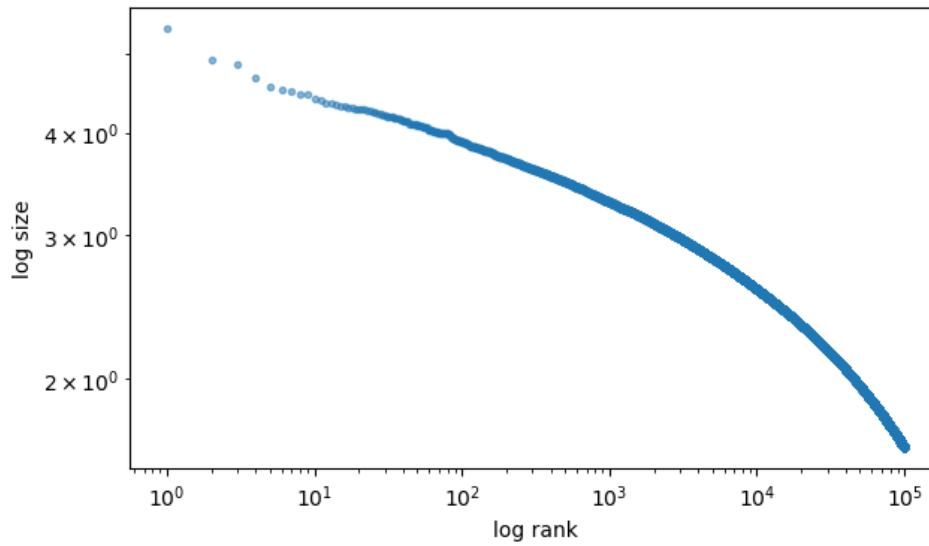
### Exercise 19.6.1

Plot the same output distribution, but this time using a rank-size plot.

If the rank-size plot is approximately linear, the data suggests a Pareto tail.

You can use QuantEcon's `rank_size` function — here's an example of usage.

```
x = jnp.abs(jax.random.normal(jax.random.PRNGKey(42), (1_000_000,)))
rank_data, size_data = qe.rank_size(x, c=0.1)
fig, ax = plt.subplots(figsize=(7,4))
ax.loglog(rank_data, size_data, "o", markersize=3.0, alpha=0.5)
ax.set_xlabel("log rank")
ax.set_ylabel("log size")
plt.show()
```

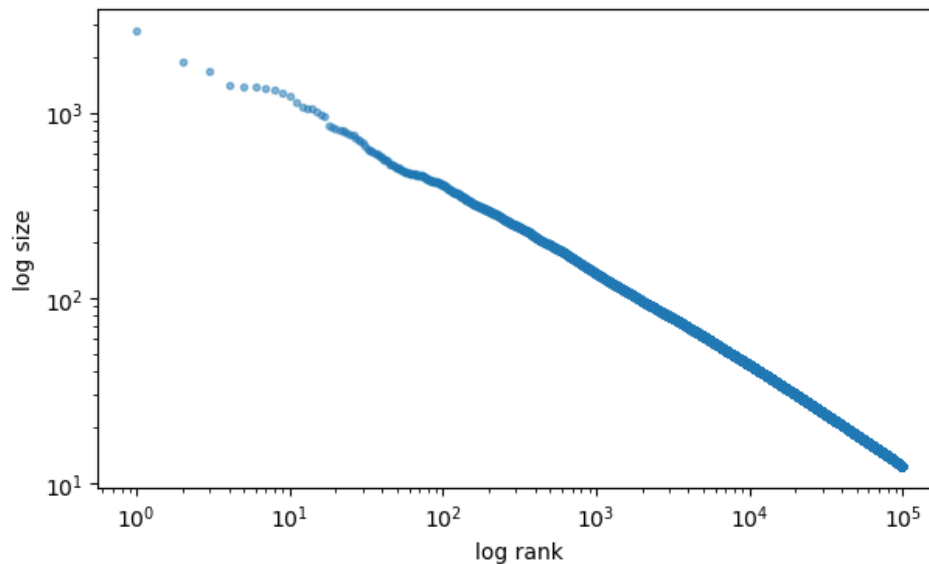


This plot is not linear — as expected, since we are using a folded normal distribution.

---

### Solution to Exercise 19.6.1

```
rank_data, size_data = qe.rank_size(output_dist, c=0.1)
fig, ax = plt.subplots(figsize=(7,4))
ax.loglog(rank_data, size_data, "o", markersize=3.0, alpha=0.5)
ax.set_xlabel("log rank")
ax.set_ylabel("log size")
plt.show()
```



This looks very linear — the model generates Pareto tails.

(In fact it's possible to prove this.)

---



**Exercise 19.6.2**

As an exercise, let's look at the fixed cost paid by incumbents each period and how it relates to the equilibrium price.

We expect that a higher fixed cost will reduce supply and hence increase the market price.

For the fixed costs, use

```
c_values = jnp.linspace(2.5, 5.0, 10)
```

**Solution to Exercise 19.6.2**

```
eq_prices = []
for i, c in enumerate(c_values):
    model = create_model(c=c)
    p_star, v_bar, psi_star, psi_sample, s, m_star = \
        compute_equilibrium_prices_and_quantities(model)
    eq_prices.append(p_star)
    print(f"Equilibrium price when c = {c:.2} is {p_star:.2}")

fig, ax = plt.subplots()
ax.plot(c_values, eq_prices, label="$p^*$")
ax.set_xlabel("fixed cost for incumbents")
ax.set_ylabel("price")
ax.legend()
plt.show()
```

Equilibrium price when c = 2.5 is 1.1

Equilibrium price when c = 2.8 is 1.2

Equilibrium price when c = 3.1 is 1.3

Equilibrium price when c = 3.3 is 1.3

Equilibrium price when c = 3.6 is 1.4

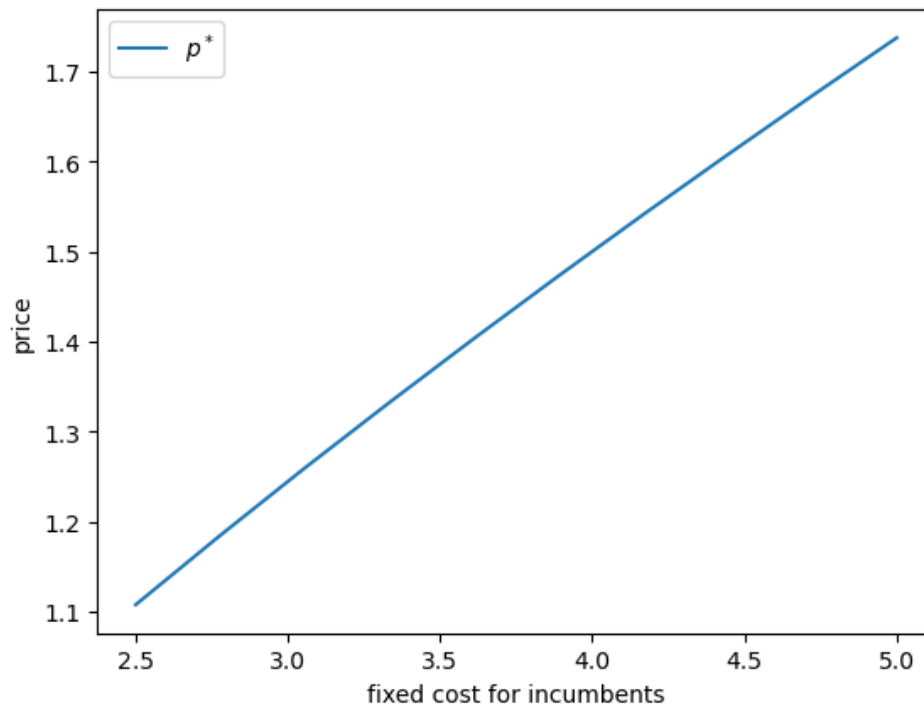
Equilibrium price when c = 3.9 is 1.5

Equilibrium price when c = 4.2 is 1.5

Equilibrium price when c = 4.4 is 1.6

Equilibrium price when c = 4.7 is 1.7

Equilibrium price when c = 5.0 is 1.7



## BIANCHI OVERBORROWING MODEL

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

This lecture provides a JAX implementation of “Overborrowing and Systemic Externalities” [[Bianchi, 2011](#)] by [Javier Bianchi](#).

In addition to what’s in Anaconda, this lecture will need the following libraries:

```
!pip install quantecon
```

We use the following imports.

```
import jax
import numba
import jax.numpy as jnp
import matplotlib.pyplot as plt
import numpy as np
import quantecon as qe
import scipy as sp
import matplotlib.pyplot as plt
import seaborn
from time import time
```

### 20.1 Description of the model

The model seeks to explain sudden stops in emerging market economies, where painful financial disruption follows a period of sustained heavy borrowing.

A representative household chooses how much to borrow on international markets and how much to consume.

The household is credit constrained, with the constraint depending on both current income and the real exchange rate.

The model shows that household overborrow because they do not internalize the effect of borrowing on the credit constraint.

This overborrowing leaves them vulnerable to bad shocks in current income.

In essence, the model works as follows

1. During good times, households borrow more and consume more
2. Increased consumption pushes up the price of nontradables and hence the real exchange rate
3. A rising real exchange rate loosens the credit constraint and encourages more borrowing
4. This leads to excessive borrowing relative to a planner.

This overborrowing leads to vulnerability vis-a-vis bad shocks.

1. When a bad shock hits, borrowing is restricted.
2. Consumption now falls, pushing down the real exchange rate.
3. This fall in the exchange rate further tightens the borrowing constraint, amplifying the shock

### 20.1.1 Decentralized equilibrium

The model contains a representative household that seeks to maximize an expected sum of discounted utility where

- the flow utility function  $u$  is CRRA, with  $u(c) = c^{1-\sigma}/(1-\sigma)$  and
- $c = (\omega c_t^{-\eta} + (1-\omega)c_n^{-\eta})^{-1/\eta}$

Here  $c_t$  ( $c_n$ ) is consumption of tradables (nontradables).

The household maximizes subject to the budget constraint

$$b' + c_t + p_n c_n = b(1+r) + y_t + p_n y_n$$

where

- $b$  is bond holdings (positive values denote assets)
- primes denote next period values
- the interest rate  $r$  is exogenous
- $p_n$  is the price of nontradables, while the price of tradables is normalized to 1
- $y_t$  and  $y_n$  are current tradable and nontradable income

The process for  $y := (y_t, y_n)$  is first-order Markov.

The household also faces the credit constraint

$$b' \geq -\kappa(y_t + p_n y_n)$$

Market clearing implies

$$c_n = y_n \quad \text{and} \quad c_t = y_t + (1+r)b - b'$$

The household takes the aggregate timepath for bonds as given by  $B' = H(B, y)$  and solves

$$V(b, B, y) = \max_{c, b'} \{u(c) + \beta \mathbb{E}_y v(b', B', y')\}$$

subject to the budget and credit constraints.

A decentralized equilibrium is a law of motion  $H$  such that the implied savings policy  $b' = b'(b, B, y)$  verifies

$$b'(B, B, y) = H(B, y) \quad \text{for all } B, y$$

### 20.1.2 Notation

Let

$$w(c_t, y_n) = \frac{c^{1-\sigma}}{1-\sigma} \quad \text{where} \quad c = [\omega c_t^{-\eta} + (1-\omega)y_n^{-\eta}]^{-1/\eta}$$

Using the market clearing conditions, we can write the household problem as

$$V(b, B, y) = \max_{b'} \{w((1+r)b + y_t - b', y_n) + \beta \mathbb{E}_y v(b', H(B, y), y')\}$$

subject to

$$-\kappa(p_n y_n + y_t) \leq b' \leq (1+r)b + y_t$$

where  $p_n$  is given by

$$p_n = ((1-\omega)/\omega)(C/y_n)^{\eta+1} \quad \text{with} \quad C := (1+r)B + y_t - H(B, y)$$

### 20.1.3 Constrained planner

The constrained planner solves

$$V(b, B, y) = \max_{c, b'} \{u(c) + \beta \mathbb{E}_y v(b', B', y')\}$$

subject to the market clearing conditions and the same constraint

$$-\kappa(y_t + p_n y_n) \leq b' \leq (1+r)b + y_t$$

although the price of nontradable is now given by

$$p_n = ((1-\omega)/\omega)(c_t/y_n)^{\eta+1} \quad \text{with} \quad c_t := (1+r)b + y_t - b'$$

We see that the planner internalizes the impact of the savings choice  $b'$  on the price of nontradables and hence the credit constraint.

## 20.2 Markov dynamics

We develop some functions for working with the VAR process

$$\ln y' = A \ln y + u' \quad (\text{prime indicates next period value})$$

where

- $y = (y_t, y_n) = (\text{tradables}, \text{nontradables})$
- $A$  is  $2 \times 2$
- $u' \sim N(0, \Omega)$
- the log function is applied pointwise

We use the following estimated values, reported on p. 12 of [Yamada \(2023\)](#).

```
A = [[0.2425, 0.3297],
      [-0.1984, 0.7576]]
Ω = [[0.0052, 0.002],
      [0.002, 0.0059]]
```

We'll store the data in  $\Omega$  using its square root:

```
C = sp.linalg.sqrtm(Ω)
A = np.array(A)
```

Here's a function to convert the VAR process to a Markov chain evolving on a rectilinear grid of points in  $\mathbb{R}^2$ .

Under the hood, this function uses the QuantEcon function `discrete_var`.

```
def discretize_income_var(A=A, C=C, grid_size=4, seed=1234):
    """
    Discretize the VAR model, returning

    y_t_nodes, a grid of y_t values
    y_n_nodes, a grid of y_n values
    Q, a Markov operator

    Let n = grid_size. The format is that Q is n x n x n x n, with

    Q[i, j, i', j'] = one step transition prob from
    (y_t_nodes[i], y_n_nodes[j]) to (y_t_nodes[i'], y_n_nodes[j'])

    """

    n = grid_size
    rng = np.random.default_rng(seed)
    mc = qe.markov.discrete_var(A, C, (n, n),
                               sim_length=1_000_000,
                               std_devs=np.sqrt(3),
                               random_state=rng)

    y_nodes, Q = np.exp(mc.state_values), mc.P
    # The array y_nodes is currently an array listing all 2 x 1 state pairs
    # (y_t, y_n), so that y_nodes[i] is the i-th such pair, while Q[l, m]
    # is the probability of transitioning from state l to state m in one step.
    # We switch the representation to the one described in the docstring.
    y_t_nodes = [y_nodes[n*i, 0] for i in range(n)]
    y_n_nodes = y_nodes[0:4, 1]
    Q = np.reshape(Q, (n, n, n, n))

    return y_t_nodes, y_n_nodes, Q
```

Here's code for sampling from the Markov chain.

```
def generate_discrete_var(A=A, C=C, grid_size=4, seed=1234,
                          ts_length=1_000_000,
                          indices=False):
    """
    Generate a time series from the discretized model, returning y_t_series and
    y_n_series. If `indices=True`, then these series are returned as grid
    indices.
    """
```

(continues on next page)

(continued from previous page)

```

n = grid_size
rng = np.random.default_rng(seed)
mc = qe.markov.discrete_var(A, C, (n, n),
                           sim_length=1_000_000,
                           std_devs=np.sqrt(3),
                           random_state=rng)

if indices:
    y_series = mc.simulate_indices(ts_length=ts_length)
    y_t_series, y_n_series = y_series % grid_size, y_series // grid_size
else:
    y_series = np.exp(mc.simulate(ts_length=ts_length))
    y_t_series, y_n_series = y_series[:, 0], y_series[:, 1]
return y_t_series, y_n_series

```

Here's code for generating the original VAR process, which can be used for testing.

```

@numba.jit
def generate_var_process(A=A, C=C, ts_length=1_000_000):
    """
    Generate the original VAR process.

    """
    y_series = np.empty((ts_length, 2))
    y_series[0, :] = np.zeros(2)
    for t in range(ts_length-1):
        y_series[t+1, :] = A @ y_series[t, :] + C @ np.random.randn(2)
    y_t_series = np.exp(y_series[:, 0])
    y_n_series = np.exp(y_series[:, 1])
    return y_t_series, y_n_series

```

Let's check some statistics for both the original and the discretized processes.

```

def corr(x, y):
    m_x, m_y = x.mean(), y.mean()
    s_xy = np.sqrt(np.sum((x - m_x)**2) * np.sum((y - m_y)**2))
    return np.sum((x - m_x) * (y - m_y)) / (s_xy)

```

```

def print_stats(y_t_series, y_n_series):
    print(f"Std dev of y_t is {y_t_series.std()}")
    print(f"Std dev of y_n is {y_n_series.std()}")
    print(f"corr(y_t, y_n) is {corr(y_t_series, y_n_series)}")
    print(f"auto_corr(y_t) is {corr(y_t_series[:-1], y_t_series[1:])}")
    print(f"auto_corr(y_n) is {corr(y_n_series[:-1], y_n_series[1:])}")
    print("\n")

```

```

print("Statistics for original process.\n")
print_stats(*generate_var_process())

```

```

Statistics for original process.

```

```
Std dev of y_t is 0.08786665450592301
Std dev of y_n is 0.10553816506894553
corr(y_t, y_n) is 0.540315803170419
auto_corr(y_t) is 0.4558023516128351
auto_corr(y_n) is 0.6648046319195241
```

```
print("Statistics for discretized process.\n")
print_stats(*generate_discrete_var())
```

```
Statistics for discretized process.
```

```
Std dev of y_t is 0.08753613842184708
Std dev of y_n is 0.10546398305554536
corr(y_t, y_n) is 0.4763676161427541
auto_corr(y_t) is 0.4013004558421011
auto_corr(y_n) is 0.5896270019961042
```

## 20.3 Overborrowing Model

In what follows

- $y = (y_t, y_n)$  is the exogenous state process

Individual states and actions are

- $c$  = consumption of tradables ( $c$  rather than  $c_t$ )
- $b$  = household savings (bond holdings)
- $bp$  = household savings decision

Aggregate quantities and prices are

- $p$  = price of nontradables ( $p$  rather than  $p_n$ )
- $B$  = aggregate savings (bond holdings)
- $C$  = aggregate consumption
- $H$  = current guess of update rule as an array of the form  $H(B, y)$

Here's code to create three tuples that store model data relevant for computation.

```
def create_overborrowing_model(
    sigma=2,                # CRRA utility parameter
    eta=(1/0.83)-1,         # Elasticity = 0.83, eta = 0.2048
    beta=0.91,              # Discount factor
    omega=0.31,             # Aggregation constant
    kappa=0.3235,           # Constraint parameter
    r=0.04,                 # Interest rate
    b_size=400,             # Bond grid size
    b_grid_min=-1.02,       # Bond grid min
    b_grid_max=-0.2         # Bond grid max (originally -0.6 to match fig)
):
    """
    Creates an instance of the overborrowing model using
```

(continues on next page)



(continued from previous page)

```

    * default parameter values from Bianchi (2011)
    * Markov dynamics from Yamada (2023)

The Markov kernel  $Q$  has the interpretation

 $Q[i, j, ip, jp]$  = one step prob of moving

 $(y_t[i], y_n[j]) \rightarrow (y_t[ip], y_n[jp])$ 

"""
# Read in Markov data and shift to JAX arrays
data = discretize_income_var()
y_t_nodes, y_n_nodes, Q = tuple(map(jnp.array, data))
# Set up grid for bond holdings
b_grid = jnp.linspace(b_grid_min, b_grid_max, b_size)
# Pack and return
parameters =  $\sigma$ ,  $\eta$ ,  $\beta$ ,  $\omega$ ,  $\kappa$ ,  $r$ 
sizes = b_size, len(y_t_nodes)
arrays = b_grid, y_t_nodes, y_n_nodes, Q
return parameters, sizes, arrays

```

Here's flow utility.

```

@jax.jit
def w(parameters, c, y_n):
    """
    Current utility when  $c_t = c$  and  $c_n = y_n$ .

     $a = [\omega c^{(-\eta)} + (1 - \omega) y_n^{(-\eta)}]^{(-1/\eta)}$ 

     $w(c, y_n) := a^{(1 - \sigma)} / (1 - \sigma)$ 

    """
     $\sigma$ ,  $\eta$ ,  $\beta$ ,  $\omega$ ,  $\kappa$ ,  $r$  = parameters
    a = ( $\omega * c^{(-\eta)} + (1 - \omega) * y_n^{(-\eta)}$ )**(-1/ $\eta$ )
    return a**((1 -  $\sigma$ ) / (1 -  $\sigma$ ))

```

We need code to generate an initial guess of  $H$ .

```

def generate_initial_H(parameters, sizes, arrays, at_constraint=False):
    """
    Compute an initial guess for  $H$ . Repeat the indices for  $b_{grid}$  over  $y_t$  and
     $y_n$  axes.

    """
    b_size, y_size = sizes
    b_indices = jnp.arange(b_size)
    O = jnp.ones((b_size, y_size, y_size), dtype=int)
    return O * jnp.reshape(b_indices, (b_size, 1, 1))

```

```
generate_initial_H = jax.jit(generate_initial_H, static_argnums=(1,))
```

We need to construct the Bellman operator for the household.

Our first function returns the (unmaximized) RHS of the Bellman equation.

```

@jax.jit
def T_generator(v, H, parameters, arrays, i_b, i_B, i_y_t, i_y_n, i_bp):
    """
    Given current state (b, B, y_t, y_n) with indices (i_b, i_B, i_y_t, i_y_n),
    compute the unmaximized right hand side (RHS) of the Bellman equation as a
    function of the next period choice bp = b', with index i_bp.
    """
    # Unpack
    σ, η, β, ω, κ, r = parameters
    b_grid, y_t_nodes, y_n_nodes, Q = arrays
    # Compute next period aggregate bonds given H
    i_Bp = H[i_B, i_y_t, i_y_n]
    # Evaluate states and actions at indices
    B, Bp, b, bp = b_grid[i_B], b_grid[i_Bp], b_grid[i_b], b_grid[i_bp]
    y_t = y_t_nodes[i_y_t]
    y_n = y_n_nodes[i_y_n]
    # Compute price of nontradables using aggregates
    C = (1 + r) * B + y_t - Bp
    p = ((1 - ω) / ω) * (C / y_n)**(η + 1)
    # Compute household flow utility
    c = (1 + r) * b + y_t - bp
    utility = w(parameters, c, y_n)
    # Compute expected value  $E_{y'} v(b', B', y')$   $Q(y, y')$ 
    EV = jnp.sum(v[i_bp, i_Bp, :, :] * Q[i_y_t, i_y_n, :, :])
    # Set up constraints
    credit_constraint_holds = - κ * (p * y_n + y_t) <= bp
    budget_constraint_holds = bp <= (1 + r) * b + y_t
    constraints_hold = jnp.logical_and(credit_constraint_holds,
                                       budget_constraint_holds)

    # Compute and return
    return jnp.where(constraints_hold, utility + β * EV, -jnp.inf)

```

Let's now vectorize and jit-compile this map.

```

# Vectorize over the control bp and all the current states
T_vec_1 = jax.vmap(T_generator,
    in_axes=(None, None, None, None, None, None, None, None, 0))
T_vec_2 = jax.vmap(T_vec_1,
    in_axes=(None, None, None, None, None, None, None, 0, None))
T_vec_3 = jax.vmap(T_vec_2,
    in_axes=(None, None, None, None, None, None, 0, None, None))
T_vec_4 = jax.vmap(T_vec_3,
    in_axes=(None, None, None, None, None, 0, None, None, None))
T_vectorized = jax.vmap(T_vec_4,
    in_axes=(None, None, None, None, 0, None, None, None, None))

```

Now we can set up the Bellman operator by maximizing over the choice variable  $b'$ .

```

def T(parameters, sizes, arrays, v, H):
    """
    Evaluate the RHS of the Bellman equation at all states and actions and then
    maximize with respect to actions.

    Return
        * Tv as an array of shape (b_size, b_size, y_size, y_size).
    """

```

(continues on next page)

(continued from previous page)

```

"""
b_size, y_size = sizes
b_grid, y_t_nodes, y_n_nodes, Q = arrays
b_indices, y_indices = jnp.arange(b_size), jnp.arange(y_size)
val = T_vectorized(v, H, parameters, arrays,
                  b_indices, b_indices, y_indices, y_indices, b_indices)
# Maximize over bp
return jnp.max(val, axis=-1)

```

```
T = jax.jit(T, static_argnums=(1,))
```

Here's a function that computes a greedy policy (best response to  $v$ ).

```

def get_greedy(parameters, sizes, arrays, v, H):
    """
    Compute the greedy policy for the household, which maximizes the right hand
    side of the Bellman equation given v and H. The greedy policy is recorded
    as an array giving the index i in b_grid such that b_grid[i] is the optimal
    choice, for every state.

    Return

    * bp_policy as an array of shape (b_size, b_size, y_size, y_size).

    """
    b_size, y_size = sizes
    b_grid, y_t_nodes, y_n_nodes, Q = arrays
    b_indices, y_indices = jnp.arange(b_size), jnp.arange(y_size)
    val = T_vectorized(v, H, parameters, arrays,
                      b_indices, b_indices, y_indices, y_indices, b_indices)
    return jnp.argmax(val, axis=-1)

```

```
get_greedy = jax.jit(get_greedy, static_argnums=(1,))
```

Here's some code for value function iteration (VFI).

```

def vfi(T, v_init, max_iter=10_000, tol=1e-5):
    """
    Use successive approximation to compute the fixed point of T, starting from
    v_init.

    """
    v = v_init

    def cond_fun(state):
        error, i, v = state
        return (error > tol) & (i < max_iter)

    def body_fun(state):
        error, i, v = state
        v_new = T(v)
        error = jnp.max(jnp.abs(v_new - v))
        return error, i+1, v_new

```

(continues on next page)

(continued from previous page)

```

error, i, v_new = jax.lax.while_loop(cond_fun, body_fun,
                                     (tol+1, 0, v))

return v_new, i

vfi = jax.jit(vfi, static_argnums=(0,))

```

This is how we update our guess of  $H$ , using the current policy  $b'$  and a damped fixed point iteration scheme.

```

def update_H(parameters, sizes, arrays, H, a):
    """
    Update guess of the aggregate update rule.

    """
    # Set up
    b_size, y_size = sizes
    b_grid, y_t_nodes, y_n_nodes, Q = arrays
    b_indices = jnp.arange(b_size)
    # Compute household response to current guess H
    v_init = jnp.ones((b_size, b_size, y_size, y_size))
    _T = lambda v: T(parameters, sizes, arrays, v, H)
    v, vfi_num_iter = vfi(_T, v_init)
    bp_policy = get_greedy(parameters, sizes, arrays, v, H)
    # Switch policy arrays to values rather than indices
    H_vals = b_grid[H]
    bp_vals = b_grid[bp_policy]
    # Update guess
    new_H_vals = a * bp_vals[b_indices, b_indices, :, :] + (1 - a) * H_vals
    # Switch back to indices
    new_H = jnp.searchsorted(b_grid, new_H_vals)
    return new_H, vfi_num_iter

```

```
update_H = jax.jit(update_H, static_argnums=(1,))
```

Now we can write code to compute an equilibrium law of motion  $H$ .

```

def compute_equilibrium(parameters, sizes, arrays,
                        a=0.5, tol=0.005, max_iter=500):
    """
    Compute the equilibrium law of motion.

    """
    H = generate_initial_H(parameters, sizes, arrays)
    error = tol + 1
    i = 0
    msgs = []
    while error > tol and i < max_iter:
        H_new, vfi_num_iter = update_H(parameters, sizes, arrays, H, a)
        msgs.append(f"VFI terminated after {vfi_num_iter} iterations.")
        error = jnp.max(jnp.abs(b_grid[H] - b_grid[H_new]))
        msgs.append(f"Updated H at iteration {i} with error {error}.")
        H = H_new
        i += 1
    if i == max_iter:
        msgs.append("Warning: Equilibrium search iteration hit upper bound.")

```

(continues on next page)

(continued from previous page)

```
print("\n".join(msgs))
return H
```

## 20.4 Planner problem

Now we switch to the planner problem.

Our first function returns the (unmaximized) RHS of the Bellman equation.

```
@jax.jit
def planner_T_generator(v, parameters, arrays, i_b, i_y_t, i_y_n, i_bp):
    """
    Given current state (b, y_t, y_n) with indices (i_b, i_y_t, i_y_n),
    compute the unmaximized right hand side (RHS) of the Bellman equation as a
    function of the next period choice bp = b'.
    """
    σ, η, β, ω, κ, r = parameters
    b_grid, y_t_nodes, y_n_nodes, Q = arrays
    y_t = y_t_nodes[i_y_t]
    y_n = y_n_nodes[i_y_n]
    b, bp = b_grid[i_b], b_grid[i_bp]
    # Compute price of nontradables using aggregates
    c = (1 + r) * b + y_t - bp
    p = ((1 - ω) / ω) * (c / y_n)**(η + 1)
    # Compute household flow utility
    utility = w(parameters, c, y_n)
    # Compute expected value (continuation)
    EV = jnp.sum(v[i_bp, :, :] * Q[i_y_t, i_y_n, :, :])
    # Set up constraints and evaluate
    credit_constraint_holds = - κ * (p * y_n + y_t) <= bp
    budget_constraint_holds = bp <= (1 + r) * b + y_t
    return jnp.where(jnp.logical_and(credit_constraint_holds,
                                     budget_constraint_holds),
                    utility + β * EV,
                    -jnp.inf)
```

```
# Vectorize over the control bp and all the current states
planner_T_vec_1 = jax.vmap(planner_T_generator,
    in_axes=(None, None, None, None, 0))
planner_T_vec_2 = jax.vmap(planner_T_vec_1,
    in_axes=(None, None, None, None, None, 0, None))
planner_T_vec_3 = jax.vmap(planner_T_vec_2,
    in_axes=(None, None, None, None, 0, None, None))
planner_T_vectorized = jax.vmap(planner_T_vec_3,
    in_axes=(None, None, None, 0, None, None, None))
```

Now we construct the Bellman operator.

```
def planner_T(parameters, sizes, arrays, v):
    b_size, y_size = sizes
    b_grid, y_t_nodes, y_n_nodes, Q = arrays
    b_indices, y_indices = jnp.arange(b_size), jnp.arange(y_size)
```

(continues on next page)

(continued from previous page)

```

# Evaluate RHS of Bellman equation at all states and actions
val = planner_T_vectorized(v, parameters, arrays,
                           b_indices, y_indices, y_indices, b_indices)
# Maximize over bp
return jnp.max(val, axis=-1)

```

```
planner_T = jax.jit(planner_T, static_argnums=(1,))
```

Here's a function that computes a greedy policy (best response to  $v$ ).

```

def planner_get_greedy(parameters, sizes, arrays, v):
    b_size, y_size = sizes
    b_grid, y_t_nodes, y_n_nodes, Q = arrays
    b_indices, y_indices = jnp.arange(b_size), jnp.arange(y_size)
    # Evaluate RHS of Bellman equation at all states and actions
    val = planner_T_vectorized(v, parameters, arrays,
                               b_indices, y_indices, y_indices, b_indices)
    # Maximize over bp
    return jnp.argmax(val, axis=-1)

```

```
planner_get_greedy = jax.jit(planner_get_greedy, static_argnums=(1,))
```

Computing the planner solution is straightforward value function iteration:

```

def compute_planner_solution(model):
    """
    Compute the constrained planner solution.

    """
    parameters, sizes, arrays = model
    b_size, y_size = sizes
    b_indices = jnp.arange(b_size)
    v_init = jnp.ones((b_size, y_size, y_size))
    _T = lambda v: planner_T(parameters, sizes, arrays, v)
    # Compute household response to current guess H
    v, vfi_num_iter = vfi(_T, v_init)
    bp_policy = planner_get_greedy(parameters, sizes, arrays, v)
    return v, bp_policy, vfi_num_iter

```

## 20.5 Numerical solution

Let's now solve the model and compare the decentralized and planner solutions

### 20.5.1 Generating solutions

Here we compute the two solutions.

```
model = create_overborrowing_model()
parameters, sizes, arrays = model
b_size, y_size = sizes
b_grid, y_t_nodes, y_n_nodes, Q = arrays
```

```
2024-08-12 04:03:55.031323: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```

```
print("Computing decentralized solution.")
start = time()
H_eq = compute_equilibrium(parameters, sizes, arrays)
diff_d_with_compile = time() - start
print(f"Computed decentralized equilibrium in {diff_d_with_compile} seconds")
```

Computing decentralized solution.

```
VFI terminated after 128 iterations.
Updated H at iteration 0 with error 0.15619052946567535.
VFI terminated after 128 iterations.
Updated H at iteration 1 with error 0.08837097883224487.
VFI terminated after 129 iterations.
Updated H at iteration 2 with error 0.0780951976776123.
VFI terminated after 128 iterations.
Updated H at iteration 3 with error 0.06987470388412476.
VFI terminated after 128 iterations.
Updated H at iteration 4 with error 0.05959904193878174.
VFI terminated after 128 iterations.
Updated H at iteration 5 with error 0.05137842893600464.
VFI terminated after 128 iterations.
Updated H at iteration 6 with error 0.04315793514251709.
VFI terminated after 128 iterations.
Updated H at iteration 7 with error 0.03904759883880615.
VFI terminated after 128 iterations.
Updated H at iteration 8 with error 0.0328822135925293.
VFI terminated after 129 iterations.
Updated H at iteration 9 with error 0.02877199649810791.
VFI terminated after 128 iterations.
Updated H at iteration 10 with error 0.024661660194396973.
VFI terminated after 128 iterations.
Updated H at iteration 11 with error 0.020551443099975586.
```

(continues on next page)

(continued from previous page)

```

VFI terminated after 128 iterations.
Updated H at iteration 12 with error 0.018496274948120117.
VFI terminated after 128 iterations.
Updated H at iteration 13 with error 0.01438605785369873.
VFI terminated after 128 iterations.
Updated H at iteration 14 with error 0.012330859899520874.
VFI terminated after 129 iterations.
Updated H at iteration 15 with error 0.010275721549987793.
VFI terminated after 128 iterations.
Updated H at iteration 16 with error 0.010275721549987793.
VFI terminated after 128 iterations.
Updated H at iteration 17 with error 0.0082206130027771.
VFI terminated after 128 iterations.
Updated H at iteration 18 with error 0.0082206130027771.
VFI terminated after 128 iterations.
Updated H at iteration 19 with error 0.0061654746532440186.
VFI terminated after 128 iterations.
Updated H at iteration 20 with error 0.006165444850921631.
VFI terminated after 128 iterations.
Updated H at iteration 21 with error 0.006165444850921631.
VFI terminated after 128 iterations.
Updated H at iteration 22 with error 0.0041103363037109375.
Computed decentralized equilibrium in 54.89348554611206 seconds

```

We run it again to get rid of compile time.

```

start = time()
H_eq = compute_equilibrium(parameters, sizes, arrays)
diff_d_without_compile = time() - start
print(f"Computed decentralized equilibrium in {diff_d_without_compile} seconds")

```

```

VFI terminated after 128 iterations.
Updated H at iteration 0 with error 0.15619052946567535.
VFI terminated after 128 iterations.
Updated H at iteration 1 with error 0.08837097883224487.
VFI terminated after 129 iterations.
Updated H at iteration 2 with error 0.0780951976776123.
VFI terminated after 128 iterations.
Updated H at iteration 3 with error 0.06987470388412476.
VFI terminated after 128 iterations.
Updated H at iteration 4 with error 0.05959904193878174.
VFI terminated after 128 iterations.
Updated H at iteration 5 with error 0.05137842893600464.
VFI terminated after 128 iterations.
Updated H at iteration 6 with error 0.04315793514251709.
VFI terminated after 128 iterations.
Updated H at iteration 7 with error 0.03904759883880615.
VFI terminated after 128 iterations.
Updated H at iteration 8 with error 0.0328822135925293.
VFI terminated after 129 iterations.
Updated H at iteration 9 with error 0.02877199649810791.
VFI terminated after 128 iterations.
Updated H at iteration 10 with error 0.024661660194396973.
VFI terminated after 128 iterations.
Updated H at iteration 11 with error 0.020551443099975586.

```

(continues on next page)



(continued from previous page)

```

VFI terminated after 128 iterations.
Updated H at iteration 12 with error 0.018496274948120117.
VFI terminated after 128 iterations.
Updated H at iteration 13 with error 0.01438605785369873.
VFI terminated after 128 iterations.
Updated H at iteration 14 with error 0.012330859899520874.
VFI terminated after 129 iterations.
Updated H at iteration 15 with error 0.010275721549987793.
VFI terminated after 128 iterations.
Updated H at iteration 16 with error 0.010275721549987793.
VFI terminated after 128 iterations.
Updated H at iteration 17 with error 0.0082206130027771.
VFI terminated after 128 iterations.
Updated H at iteration 18 with error 0.0082206130027771.
VFI terminated after 128 iterations.
Updated H at iteration 19 with error 0.0061654746532440186.
VFI terminated after 128 iterations.
Updated H at iteration 20 with error 0.006165444850921631.
VFI terminated after 128 iterations.
Updated H at iteration 21 with error 0.006165444850921631.
VFI terminated after 128 iterations.
Updated H at iteration 22 with error 0.0041103363037109375.
Computed decentralized equilibrium in 54.168943643569946 seconds

```

```

print("Computing planner's solution.")
start = time()
planner_v, H_plan, vfi_num_iter = compute_planner_solution(model)
diff_p_with_compile = time() - start
print(f"Computed planner's equilibrium in {diff_p_with_compile} seconds")

```

```
Computing planner's solution.
```

```
Computed planner's equilibrium in 0.8798391819000244 seconds
```

We run it again to eliminate compile time.

```

start = time()
planner_v, H_plan, vfi_num_iter = compute_planner_solution(model)
diff_p_without_compile = time() - start
print(f"Computed planner's equilibrium in {diff_p_without_compile} seconds")

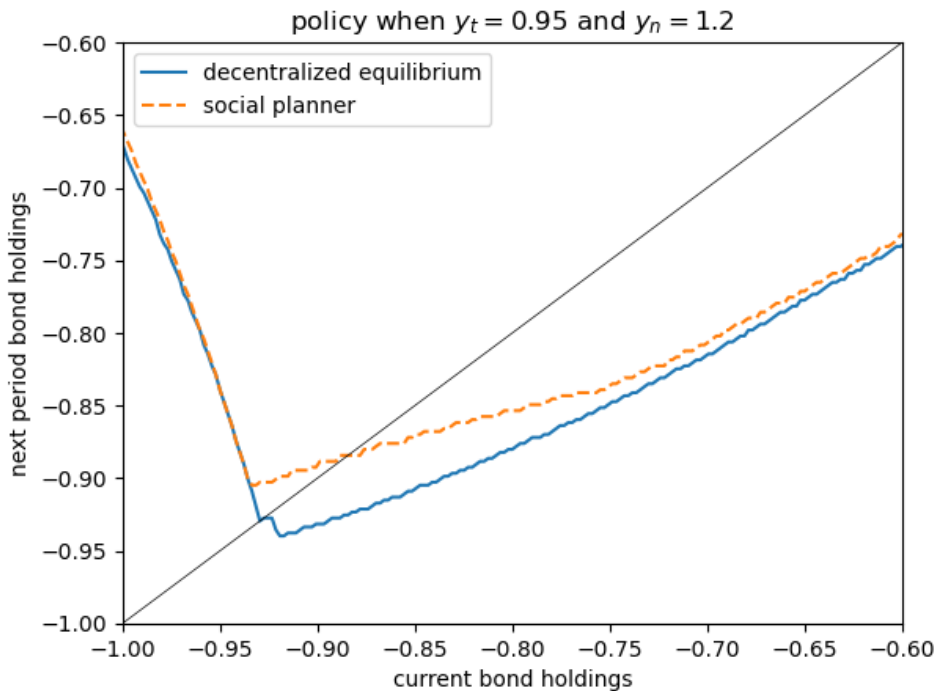
```

```
Computed planner's equilibrium in 0.22881293296813965 seconds
```

## 20.5.2 Policy plots

We produce a policy plot that is similar to Figure 1 in [Bianchi, 2011].

```
i, j = 1, 3
y_t, y_n = y_t_nodes[i], y_n_nodes[j]
fig, ax = plt.subplots()
ax.plot(b_grid, b_grid[H_eq[:, i, j]], label='decentralized equilibrium')
ax.plot(b_grid, b_grid[H_plan[:, i, j]], ls='--', label='social planner')
ax.plot(b_grid, b_grid, color='black', lw=0.5)
ax.set_ylim((-1.0, -0.6))
ax.set_xlim((-1.0, -0.6))
ax.set_xlabel("current bond holdings")
ax.set_ylabel("next period bond holdings")
ax.set_title(f"policy when $y_t = {y_t:.2}$ and $y_n = {y_n:.2}$")
ax.legend()
plt.show()
```



The match is not exact because we use a different estimate of the Markov dynamics for income. Nonetheless, it is qualitatively similar.

## 20.6 Exercise

### Exercise 20.6.1

Your task is to examine the ergodic distribution of borrowing in the decentralized and planner equilibria.

We recommend that you use simulation and a kernel density estimate.

Here's a function for generating the borrowing sequence.

We use Numba because we want to compile a long for loop.

```
@numba.jit
def generate_borrowing_sequence(H, y_t_series, y_n_series):
    """
    Generate the borrowing sequence  $B' = H(B, y_t, y_n)$ .

    *  $H$  is a policy array
    *  $y_t\_series$  and  $y_n\_series$  are simulated time paths

    Both  $y_t\_series$  and  $y_n\_series$  are stored as indices rather than values.

    """
    B = np.empty_like(y_t_series)
    B[0] = 0
    for t in range(len(y_t_series)-1):
        B[t+1] = H[B[t], y_t_series[t], y_n_series[t]]
    return B
```

Note that you will have to convert JAX arrays into NumPy arrays if you want to use this function.

From here you will need to

- generate a time path for income  $y = (y_t, y_n)$  using one of the functions provided above.
- use the function `generate_borrowing_sequence` plus `H_eq` and `H_plan` to calculate bond holdings for the planner and the decentralized equilibrium
- produce a kernel density plot for each of these data sets

If you are successful, your plot should look something like Fig 2 of [Bianchi, 2011] — although not exactly the same, due to the alternative specification of the Markov process.

To generate a kernel density plot, we recommend that you use `kdeplot` from the package `seaborn`, which is included in Anaconda.

### Solution to Exercise 20.6.1

```
sim_length = 100_000
y_t_series, y_n_series = generate_discrete_var(ts_length=sim_length,
                                              indices=True)
```

We convert JAX arrays to NumPy arrays in order to use Numba.

```
y_t_series, y_n_series, H_eq, H_plan = \
    [np.array(v) for v in (y_t_series, y_n_series, H_eq, H_plan)]
```

Now let's compute borrowing for the decentralized equilibrium and the planner.

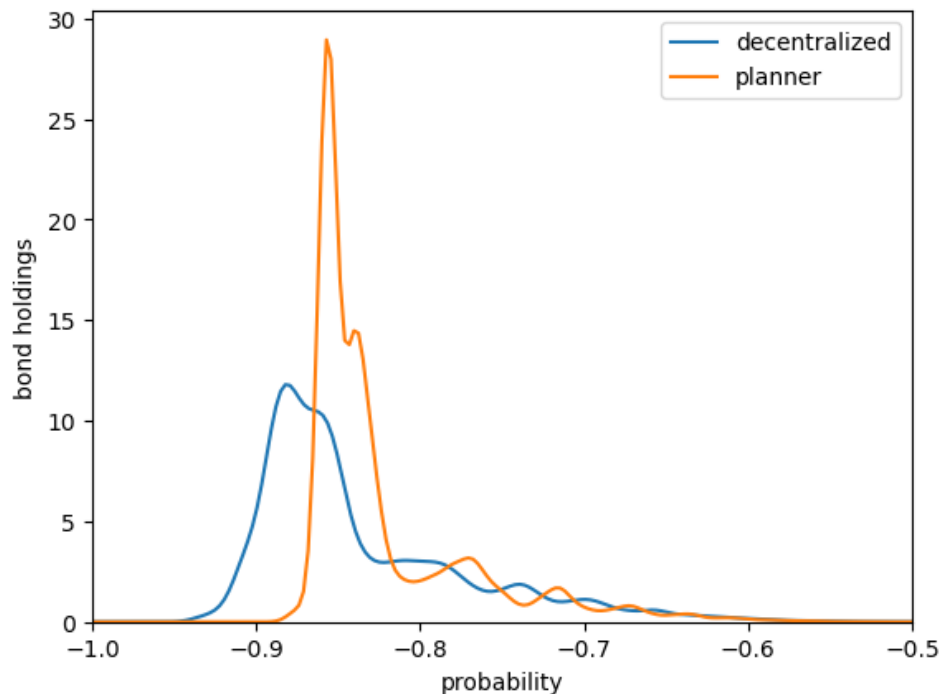
```
B_eq = generate_borrowing_sequence(H_eq, y_t_series, y_n_series)
eq_b_sequence = b_grid[B_eq]
B_plan = generate_borrowing_sequence(H_plan, y_t_series, y_n_series)
plan_b_sequence = b_grid[B_plan]
```

We suppress some annoying warnings produced by the current version of seaborn

```
import warnings
warnings.simplefilter(action='ignore', category=FutureWarning)
```

Now let's plot the distributions using a kernel density estimator.

```
fig, ax = plt.subplots()
seaborn.kdeplot(eq_b_sequence, ax=ax, label='decentralized')
seaborn.kdeplot(plan_b_sequence, ax=ax, label='planner')
ax.legend()
ax.set_xlim((-1, -0.5))
ax.set_xlabel("probability")
ax.set_ylabel("bond holdings")
plt.show()
```



This corresponds to Figure 2 in [\[Bianchi, 2011\]](#).

Again, the match is not exact but it is qualitatively similar.

Asset holding has a longer left hand tail under the decentralized equilibrium, leaving the economy more vulnerable to bad shocks.

---

**Part VI**

**Data and Empirics**



## MAXIMUM LIKELIHOOD ESTIMATION

---

### GPU

This lecture was built using a machine with JAX installed and access to a GPU.

To run this lecture on [Google Colab](#), click on the “play” icon top right, select Colab, and set the runtime environment to include a GPU.

To run this lecture on your own machine, you need to install [Google JAX](#).

---

## 21.1 Overview

This lecture is the extended JAX implementation of [this section](#) of [this lecture](#).

Please refer that lecture for all background and notation.

Here we will exploit the automatic differentiation capabilities of JAX rather than calculating derivatives by hand.

We’ll require the following imports:

```
import matplotlib.pyplot as plt
from collections import namedtuple
import jax.numpy as jnp
import jax
from statsmodels.api import Poisson
```

Let’s check the GPU we are running

```
!nvidia-smi
```

```
Mon Aug 12 04:00:10 2024
+-----+
| NVIDIA-SMI 535.54.03                  Driver Version: 535.54.03      CUDA Version: 12.
| 5                                    |
+-----+-----+-----+-----+
| GPU   Name                               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC   |
| Fan  Temp  Perf    Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
+-----+-----+-----+-----+
| Compute M. |
```

(continues on next page)





```
@jax.jit
def logL( $\beta$ ):
    return -( $\beta$  - 10) ** 2 - 10
```

To find the value of  $\frac{d \log \mathcal{L}(\beta)}{d\beta}$ , we can use `jax.grad` which auto-differentiates the given function.

We further use `jax.vmap` which vectorizes the given function i.e. the function acting upon scalar inputs can now be used with vector inputs.

```
dlogL = jax.vmap(jax.grad(logL))
```

```
 $\beta$  = jnp.linspace(1, 20)

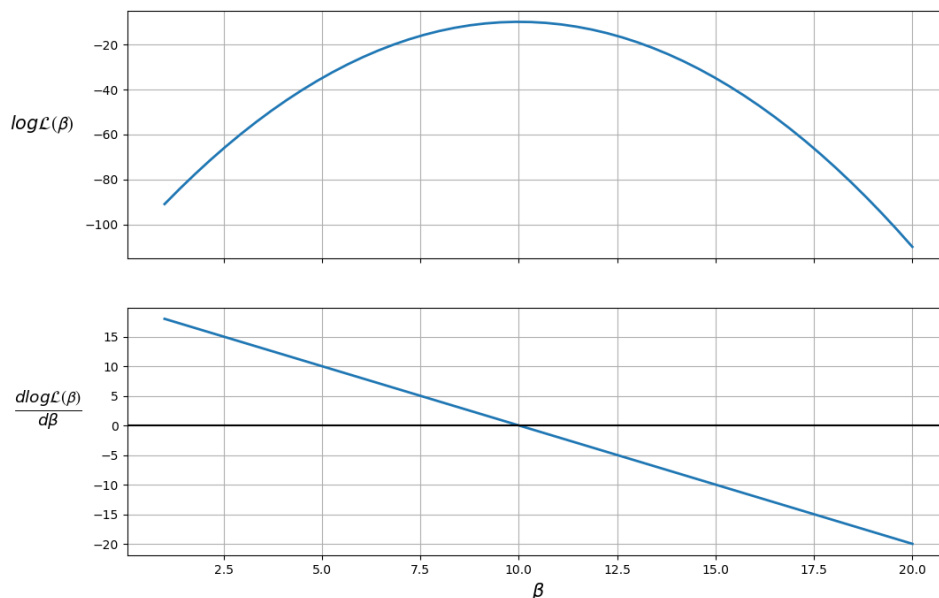
fig, (ax1, ax2) = plt.subplots(2, sharex=True, figsize=(12, 8))

ax1.plot( $\beta$ , logL( $\beta$ ), lw=2)
ax2.plot( $\beta$ , dlogL( $\beta$ ), lw=2)

ax1.set_ylabel(r'$\log \mathcal{L}(\beta)$',
               rotation=0,
               labelpad=35,
               fontsize=15)
ax2.set_ylabel(r'$\frac{d \log \mathcal{L}(\beta)}{d \beta}$ ',
               rotation=0,
               labelpad=35,
               fontsize=19)

ax2.set_xlabel(r'$\beta$', fontsize=15)
ax1.grid(), ax2.grid()
plt.axhline(c='black')
plt.show()
```

```
2024-08-12 04:00:11.155290: W external/xla/xla/service/gpu/nvptx_compiler.cc:836]
↳The NVIDIA driver's CUDA version is 12.5 which is older than the PTX compiler
↳version (12.6.20). Because the driver is older than the PTX compiler version,
↳XLA is disabling parallel compilation, which may slow down compilation. You
↳should update your NVIDIA driver or use the NVIDIA-provided CUDA forward
↳compatibility packages.
```



The plot shows that the maximum likelihood value (the top plot) occurs when  $\frac{d \log \mathcal{L}(\beta)}{d\beta} = 0$  (the bottom plot).

Therefore, the likelihood is maximized when  $\beta = 10$ .

We can also ensure that this value is a *maximum* (as opposed to a minimum) by checking that the second derivative (slope of the bottom plot) is negative.

The Newton-Raphson algorithm finds a point where the first derivative is 0.

To use the algorithm, we take an initial guess at the maximum value,  $\beta_0$  (the OLS parameter estimates might be a reasonable guess).

Then we use the updating rule involving gradient information to iterate the algorithm until the error is sufficiently small or the algorithm reaches the maximum number of iterations.

Please refer to [this section](#) for the detailed algorithm.

## 21.2.2 A Poisson model

Let's have a go at implementing the Newton-Raphson algorithm to calculate the maximum likelihood estimations of a Poisson regression.

The Poisson regression has a joint pmf:

$$f(y_1, y_2, \dots, y_n \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n; \beta) = \prod_{i=1}^n \frac{\mu_i^{y_i}}{y_i!} e^{-\mu_i}$$

$$\text{where } \mu_i = \exp(\mathbf{x}_i' \beta) = \exp(\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik})$$

We create a `namedtuple` to store the observed values

```
RegressionModel = namedtuple('RegressionModel', ['X', 'y'])

def create_regression_model(X, y):
    n, k = X.shape
    # Reshape y as a n_by_1 column vector
    y = y.reshape(n, 1)
```

(continues on next page)

(continued from previous page)

```
X, y = jax.device_put((X, y))
return RegressionModel(X=X, y=y)
```

The log likelihood function of the Poisson regression is

$$\max_{\beta} \left( \sum_{i=1}^n y_i \log \mu_i - \sum_{i=1}^n \mu_i - \sum_{i=1}^n \log y_i! \right)$$

The full derivation can be found [here](#).

The log likelihood function involves factorial, but JAX doesn't have a readily available implementation to compute factorial directly.

In order to compute the factorial efficiently such that we can JIT it, we use

$$n! = e^{\log(\Gamma(n+1))}$$

since `jax.lax.lgamma` and `jax.lax.exp` are available.

The following function `jax_factorial` computes the factorial using this idea.

Let's define this function in Python

```
@jax.jit
def _factorial(n):
    return jax.lax.exp(jax.lax.lgamma(n + 1.0)).astype(int)

jax_factorial = jax.vmap(_factorial)
```

Now we can define the log likelihood function in Python

```
@jax.jit
def poisson_logL(beta, model):
    y = model.y
    mu = jnp.exp(model.X @ beta)
    return jnp.sum(model.y * jnp.log(mu) - mu - jnp.log(jax_factorial(y)))
```

To find the gradient of the `poisson_logL`, we again use `jax.grad`.

According to [the documentation](#),

- `jax.jacfwd` uses forward-mode automatic differentiation, which is more efficient for “tall” Jacobian matrices, while
- `jax.jacrev` uses reverse-mode, which is more efficient for “wide” Jacobian matrices.

(The documentation also states that when matrices that are near-square, `jax.jacfwd` probably has an edge over `jax.jacrev`.)

Therefore, to find the Hessian, we can directly use `jax.jacfwd`.

```
G_poisson_logL = jax.grad(poisson_logL)
H_poisson_logL = jax.jacfwd(G_poisson_logL)
```

Our function `newton_raphson` will take a `RegressionModel` object that has an initial guess of the parameter vector  $\beta_0$ .

The algorithm will update the parameter vector according to the updating rule, and recalculate the gradient and Hessian matrices at the new parameter estimates.

```

def newton_raphson(model,  $\beta$ , tol=1e-3, max_iter=100, display=True):

    i = 0
    error = 100  # Initial error value

    # Print header of output
    if display:
        header = f'{"Iteration_k":<13}{"Log-likelihood":<16}{" $\theta$ ":<60}'
        print(header)
        print("-" * len(header))

    # While loop runs while any value in error is greater
    # than the tolerance until max iterations are reached
    while jnp.any(error > tol) and i < max_iter:
        H, G = jnp.squeeze(H_poisson_logL( $\beta$ , model)), G_poisson_logL( $\beta$ , model)
         $\beta_{\text{new}}$  =  $\beta$  - (jnp.dot(jnp.linalg.inv(H), G))
        error = jnp.abs( $\beta_{\text{new}}$  -  $\beta$ )
         $\beta$  =  $\beta_{\text{new}}$ 

        if display:
             $\beta_{\text{list}}$  = [f'{t:.3}' for t in list( $\beta$ .flatten())]
            update = f'{"i":<13}{"poisson_logL( $\beta$ , model):<16.8}{" $\beta_{\text{list}}$ "}'
            print(update)

        i += 1

    print(f'Number of iterations: {i}')
    print(f' $\beta_{\text{hat}}$  = { $\beta$ .flatten()}')

    return  $\beta$ 

```

Let's try out our algorithm with a small dataset of 5 observations and 3 variables in **X**.

```

X = jnp.array([[1, 2, 5],
               [1, 1, 3],
               [1, 4, 2],
               [1, 5, 2],
               [1, 3, 1]])

y = jnp.array([1, 0, 1, 1, 0])

# Take a guess at initial  $\beta$ s
init_ $\beta$  = jnp.array([0.1, 0.1, 0.1]).reshape(X.shape[1], 1)

# Create an object with Poisson model values
poi = create_regression_model(X, y)

# Use newton_raphson to find the MLE
 $\beta_{\text{hat}}$  = newton_raphson(poi, init_ $\beta$ , display=True)

```

```

Iteration_k  Log-likelihood   $\theta$ 
-----

```

```

↪-----

```

```

0           -4.3447622      ['-1.49', '0.265', '0.244']

```

(continues on next page)

(continued from previous page)

```

1          -3.5742413      ['-3.38', '0.528', '0.474']
2          -3.3999526      ['-5.06', '0.782', '0.702']
3          -3.3788646      ['-5.92', '0.909', '0.82']
4          -3.3783559      ['-6.07', '0.933', '0.843']
5          -3.3783555      ['-6.08', '0.933', '0.843']
6          -3.3783555      ['-6.08', '0.933', '0.843']
Number of iterations: 7
β_hat = [-6.07848573  0.9334028  0.84329677]

```

As this was a simple model with few observations, the algorithm achieved convergence in only 7 iterations.

The gradient vector should be close to 0 at  $\hat{\beta}$

```
G_poisson_logL(β_hat, poi)
```

```

Array([[ -2.55892529e-13],
       [-6.50313137e-13],
       [-5.01695907e-13]], dtype=float64)

```

## 21.3 MLE with statsmodels

We'll use the Poisson regression model in `statsmodels` to verify the results obtained using JAX.

`statsmodels` uses the same algorithm as above to find the maximum likelihood estimates.

Now, as `statsmodels` accepts only NumPy arrays, we can use the `__array__` method of JAX arrays to convert it to NumPy arrays.

```

X_numpy = X.__array__()
y_numpy = y.__array__()

```

```

stats_poisson = Poisson(y_numpy, X_numpy).fit()
print(stats_poisson.summary())

```

```

Optimization terminated successfully.
    Current function value: 0.675671
    Iterations 7

                        Poisson Regression Results
=====
Dep. Variable:          y      No. Observations:          5
Model:                Poisson  Df Residuals:              2
Method:                  MLE    Df Model:                2
Date:                  Mon, 12 Aug 2024    Pseudo R-squ.:      0.2546
Time:                  04:00:13    Log-Likelihood:      -3.3784
converged:              True    LL-Null:             -4.5325
Covariance Type:        nonrobust    LLR p-value:         0.3153
=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
const         -6.0785         5.279        -1.151      0.250      -16.425         4.268
x1              0.9334         0.829         1.126      0.260       -0.691         2.558

```

(continues on next page)

(continued from previous page)

x2	0.8433	0.798	1.057	0.291	-0.720	2.407
=====						

The benefit of writing our own procedure, relative to statsmodels is that

- we can exploit the power of the GPU and
- we learn the underlying methodology, which can be extended to complex situations where no existing routines are available.

### Exercise 21.3.1

We define a quadratic model for a single explanatory variable by

$$\log(\lambda_t) = \beta_0 + \beta_1 x_t + \beta_2 x_t^2$$

We calculate the mean on the original scale instead of the log scale by exponentiating both sides of the above equation, which gives

$$\lambda_t = \exp(\beta_0 + \beta_1 x_t + \beta_2 x_t^2) \quad (21.1)$$

Simulate the values of  $x_t$  by sampling from a normal distribution and  $\lambda_t$  by using (21.1) and the following constants:

$$\beta_0 = -2.5, \quad \beta_1 = 0.25, \quad \beta_2 = 0.5$$

Try to obtain the approximate values of  $\beta_0, \beta_1, \beta_2$ , by simulating a Poisson Regression Model such that

$$y_t \sim \text{Poisson}(\lambda_t) \quad \text{for all } t.$$

Using our `newton_raphson` function on the data set  $X = [1, x_t, x_t^2]$  and  $y$ , obtain the maximum likelihood estimates of  $\beta_0, \beta_1, \beta_2$ .

With a sufficient large sample size, you should approximately recover the true values of these parameters.

### Solution to Exercise 21.3.1

Let's start by defining "true" parameter values.

```
β_0 = -2.5
β_1 = 0.25
β_2 = 0.5
```

To simulate the model, we sample 500,000 values of  $x_t$  from the standard normal distribution.

```
seed = 32
shape = (500_000, 1)
key = jax.random.PRNGKey(seed)
x = jax.random.normal(key, shape)
```

We compute  $\lambda$  using (21.1)

```
λ = jnp.exp(β_0 + β_1 * x + β_2 * x**2)
```

Let's define  $y_t$  by sampling from a Poisson distribution with mean as  $\lambda_t$ .

```
y = jax.random.poisson(key,  $\lambda$ , shape)
```

Now let's try to recover the true parameter values using the Newton-Raphson method described above.

```
X = jnp.hstack((jnp.ones(shape), x, x**2))

# Take a guess at initial  $\beta$ s
init_β = jnp.array([0.1, 0.1, 0.1]).reshape(X.shape[1], 1)

# Create an object with Poisson model values
poi = create_regression_model(X, y)

# Use newton_raphson to find the MLE
β_hat = newton_raphson(poi, init_β, tol=1e-5, display=True)
```

```
Iteration_k  Log-likelihood   $\theta$ 
```

```
-----
```

```
0          -4.5444745e+07  ['-1.49', '0.312', '0.794']
1          -1.6303734e+07  ['-2.42', '0.311', '0.79']
2          -5689832.6      ['-3.22', '0.31', '0.78']
3          -1869457.7      ['-3.73', '0.306', '0.756']
4          -510284.64      ['-3.71', '0.297', '0.705']
5          -28066.381      ['-3.2', '0.283', '0.63']
6           127470.33      ['-2.77', '0.268', '0.563']
7           163044.64      ['-2.57', '0.257', '0.52']
8           166608.12      ['-2.51', '0.252', '0.503']
9           166666.4       ['-2.5', '0.251', '0.5']
10          166666.42      ['-2.5', '0.251', '0.5']
11          166666.42      ['-2.5', '0.251', '0.5']
Number of iterations: 12
β_hat = [-2.50016027  0.25079345  0.50008394]
```

The maximum likelihood estimates are similar to the true parameter values.





## **Part VII**

### **Other**



## TROUBLESHOOTING

This page is for readers experiencing errors when running the code from the lectures.

### 22.1 Fixing Your Local Environment

The basic assumption of the lectures is that code in a lecture should execute whenever

1. it is executed in a Jupyter notebook and
2. the notebook is running on a machine with the latest version of Anaconda Python.

You have installed Anaconda, haven't you, following the instructions in [this lecture](#)?

Assuming that you have, the most common source of problems for our readers is that their Anaconda distribution is not up to date.

[Here's a useful article](#) on how to update Anaconda.

Another option is to simply remove Anaconda and reinstall.

You also need to keep the external code libraries, such as [QuantEcon.py](#) up to date.

For this task you can either

- use `conda install -y quantecon` on the command line, or
- execute `!conda install -y quantecon` within a Jupyter notebook.

If your local environment is still not working you can do two things.

First, you can use a remote machine instead, by clicking on the Launch Notebook icon available for each lecture



Second, you can report an issue, so we can try to fix your local set up.

We like getting feedback on the lectures so please don't hesitate to get in touch.

## 22.2 Reporting an Issue

One way to give feedback is to raise an issue through our [issue tracker](#).

Please be as specific as possible. Tell us where the problem is and as much detail about your local set up as you can provide.

Another feedback option is to use our [discourse forum](#).

Finally, you can provide direct feedback to [contact@quantecon.org](mailto:contact@quantecon.org)

---

CHAPTER  
**TWENTYTHREE**

---

**REFERENCES**



## EXECUTION STATISTICS

This table contains the latest execution statistics.

Document	Modified	Method	Run Time (s)	Status
<a href="#">aiyagari_jax</a>	2024-08-12 03:21	cache	68.16	✓
<a href="#">arellano</a>	2024-08-12 03:21	cache	20.5	✓
<a href="#">autodiff</a>	2024-08-12 03:21	cache	9.53	✓
<a href="#">hopenhayn</a>	2024-08-12 03:24	cache	130.79	✓
<a href="#">ifp_egm</a>	2024-08-12 03:26	cache	119.51	✓
<a href="#">intro</a>	2024-08-12 03:26	cache	1.09	✓
<a href="#">inventory_dynamics</a>	2024-08-12 03:26	cache	39.45	✓
<a href="#">inventory_ssd</a>	2024-08-12 03:58	cache	1920.47	✓
<a href="#">jax_intro</a>	2024-08-12 03:59	cache	27.78	✓
<a href="#">job_search</a>	2024-08-12 03:59	cache	9.53	✓
<a href="#">kesten_processes</a>	2024-08-12 03:59	cache	13.41	✓
<a href="#">lucas_model</a>	2024-08-12 03:59	cache	15.49	✓
<a href="#">markov_asset</a>	2024-08-12 04:00	cache	11.48	✓
<a href="#">mle</a>	2024-08-12 04:00	cache	14.49	✓
<a href="#">newtons_method</a>	2024-08-12 04:02	cache	125.57	✓
<a href="#">opt_invest</a>	2024-08-12 04:02	cache	22.94	✓
<a href="#">opt_savings_1</a>	2024-08-12 04:03	cache	31.92	✓
<a href="#">opt_savings_2</a>	2024-08-12 04:03	cache	20.99	✓
<a href="#">overborrowing</a>	2024-08-12 04:05	cache	123.48	✓
<a href="#">short_path</a>	2024-08-12 04:05	cache	3.02	✓
<a href="#">status</a>	2024-08-12 04:05	cache	2.11	✓
<a href="#">troubleshooting</a>	2024-08-12 03:26	cache	1.09	✓
<a href="#">wealth_dynamics</a>	2024-08-12 04:09	cache	216.57	✓
<a href="#">zreferences</a>	2024-08-12 03:26	cache	1.09	✓

These lectures are built on `linux` instances through `github actions` that has access to a `gpu`. These lectures make use of the `nvidia T4` card.

You can check the backend used by JAX using:

```
import jax
# Check if JAX is using GPU
print(f"JAX backend: {jax.devices()[0].platform}")
```

```
JAX backend: gpu
```

and the hardware we are running on:

```
!nvidia-smi
```

```
/opt/conda/envs/quantecon/lib/python3.11/pty.py:89: RuntimeWarning: os.fork() was
called. os.fork() is incompatible with multithreaded code, and JAX is
multithreaded, so this will likely lead to a deadlock.
pid, fd = os.forkpty()
```

```
Mon Aug 12 04:05:53 2024
+-----+
| NVIDIA-SMI 535.54.03              Driver Version: 535.54.03      CUDA Version: 12.
5      |
+-----+-----+-----+-----+
| GPU   Name                     Persistence-M | Bus-Id        Disp.A | Volatile Uncorr.
| ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap |      Memory-Usage | GPU-Util  |
| Compute M. |
|
| MIG M. |
+-----+-----+-----+-----+
|   0   Tesla T4                       On      | 00000001:00:00.0 Off  |
| Off |
| N/A   62C    P0               37W / 70W | 107MiB / 16384MiB |      2%
| Default |
|
| N/A |
+-----+-----+-----+-----+
+-----+
| Processes:
|
| GPU   GI    CI        PID    Type    Process name                        GPU
| Memory |
| ID     ID                               Usage
+-----+-----+-----+-----+
+-----+
```



## BIBLIOGRAPHY

- [Are08] Cristina Arellano. Default risk and income fluctuations in emerging economies. *The American Economic Review*, pages 690–712, 2008.
- [Bia11] Javier Bianchi. Overborrowing and systemic externalities in the business cycle. *American Economic Review*, 101(7):3400–3426, December 2011. URL: <https://www.aeaweb.org/articles?id=10.1257/aer.101.7.3400>, doi:10.1257/aer.101.7.3400.
- [Luc78] Robert E Lucas, Jr. Asset prices in an exchange economy. *Econometrica: Journal of the Econometric Society*, 46(6):1429–1445, 1978.



## INDEX

### F

Fixed Point Theory, [101](#)

### K

Kesten processes  
    heavy tails, [68](#)

### L

Linear State Space Models, [67](#)

Lucas Model, [98](#)

- Assets, [98](#)
- Computation, [102](#)
- Consumers, [98](#)
- Dynamic Program, [99](#)
- Equilibrium Constraints, [100](#)
- Equilibrium Price Function, [100](#)
- Pricing, [99](#)
- Solving, [101](#)

### M

Markov process, inventory, [55](#)

Models

- Lucas Asset Pricing, [97](#)