

# Lecture 9: Simulation and Inference

Zhuotong Liu

University of Edinburgh

Autumn 2025

# Using our models

- ▶ There are three main reasons why a solved model is useful:
  1. **Understanding the mechanisms:** sometimes, plotting the solutions or simulating data can be illuminating about how a model works
  2. **Estimating Parameters:** We want to choose parameters of the model in a sensible way so that our model matches the data
  3. **Counterfactual Exercises:** What happens if we change something in the model? What does it predict? How does it change welfare? Other metrics we care about?
- ▶ It turns out that to do all of this, we need to be able to generate fake (simulated) data from our models

# Section 1

## Simulating Data from Models

## Simulating Data: Deterministic Case

- ▶ Consider the non-stochastic neoclassical growth model we saw last week:

$$\begin{aligned} v(k) &= \max_{c, k'} u(c) + \beta v(k') \\ \text{s.t.} \quad & c + k' \leq Ak^\alpha + (1 - \delta)k \end{aligned} \tag{1}$$

- ▶ A solution to this model will give us some representation of  $v(k)$  as well as policy functions  $g_c(k)$  and  $g_k(k)$ , where

$$v(k) = u(g_c(k)) + \beta v(g_k(k))$$

- ▶ Suppose we start at an arbitrary  $k_0$ .
- ▶ Since the model is deterministic, generating “data” from the model is just a matter of stepping our value  $k_0$  forward one period at a time using our policy functions

$$k_t = g_k(k_{t-1})$$

$$c_t = g_c(k_t)$$

## Simulating Data: Deterministic Case

- ▶ Consider the non-stochastic neoclassical growth model we saw last week:

$$\begin{aligned} v(k) &= \max_{c, k'} u(c) + \beta v(k') \\ \text{s.t.} \quad & c + k' \leq Ak^\alpha + (1 - \delta)k \end{aligned} \tag{1}$$

- ▶ A solution to this model will give us some representation of  $v(k)$  as well as policy functions  $g_c(k)$  and  $g_k(k)$ , where

$$v(k) = u(g_c(k)) + \beta v(g_k(k))$$

- ▶ Suppose we start at an arbitrary  $k_0$ .
- ▶ Since the model is deterministic, generating “data” from the model is just a matter of stepping our value  $k_0$  forward one period at a time using our policy functions

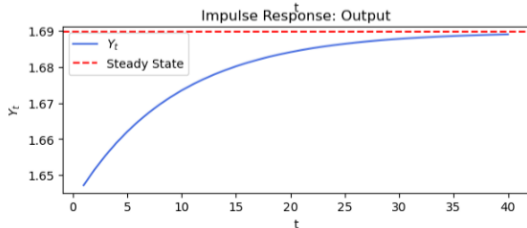
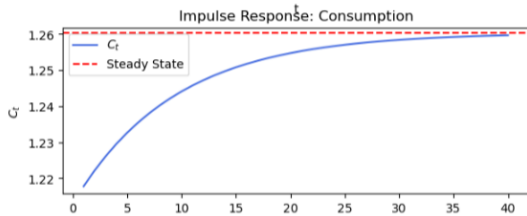
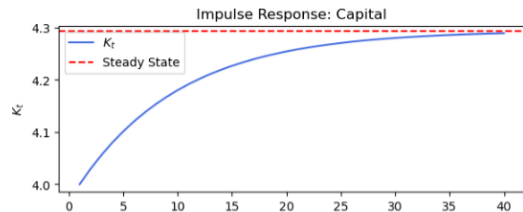
$$k_t = g_k(k_{t-1})$$

$$c_t = g_c(k_t)$$

# Stepping the Model Forward

```
def simulate(p, policy, kgrid, k0, T=40):
    A, alpha, delta, beta = p
    C = np.zeros(T)
    K = np.zeros(T)
    Y = np.zeros(T)

    K[0] = k0
    for t in range(T):
        Y[t] = A * K[t]**alpha
        z = Y[t] + (1 - delta) * K[t]      # Available resources
        k_next = policy(K[t])             # Next-period capital
        k_next = np.clip(k_next, kgrid[0], kgrid[-1])
        C[t] = z - k_next                 # Consumption
        if t < T-1:
            K[t+1] = k_next
    return C, K, Y
```



## Non-deterministic case: Simulate data

- ▶ Suppose that we return to the case of our stochastic neoclassical growth model:

$$\begin{aligned} v(k, A) = \max_{c, k'} & \quad u(c) + \beta \mathbb{E} [v(k', A') | A] \\ \text{s.t.} & \quad c + k' \leq Ak^\alpha + (1 - \delta)k \\ & \quad \log(A') = \rho \log(A) + \epsilon \\ & \quad \epsilon \sim N(0, \sigma) \end{aligned} \tag{2}$$

- ▶ How can we simulate data from this model?

1. Solve the model, and recover the policy rules  $g_c(k, A)$  and  $g_k(k, A)$
2. Start from  $(k_0, A_0)$
3. Draw a random sequence of  $A_t$  that satisfy our process for  $A$
4. For each  $t > 0$ , set

$$\begin{aligned} k_t &= g_k(k_{t-1}, A_{t-1}) \\ c_t &= g_c(k_t, A_t) \end{aligned}$$

## Non-deterministic case: Simulate data

- ▶ Suppose that we return to the case of our stochastic neoclassical growth model:

$$\begin{aligned} v(k, A) = \max_{c, k'} & \quad u(c) + \beta \mathbb{E} [v(k', A') | A] \\ \text{s.t.} & \quad c + k' \leq Ak^\alpha + (1 - \delta)k \\ & \quad \log(A') = \rho \log(A) + \epsilon \\ & \quad \epsilon \sim N(0, \sigma) \end{aligned} \tag{2}$$

- ▶ How can we simulate data from this model?

1. Solve the model, and recover the policy rules  $g_c(k, A)$  and  $g_k(k, A)$
2. Start from  $(k_0, A_0)$
3. Draw a random sequence of  $A_t$  that satisfy our process for  $A$
4. For each  $t > 0$ , set

$$\begin{aligned} k_t &= g_k(k_{t-1}, A_{t-1}) \\ c_t &= g_c(k_t, A_t) \end{aligned}$$

## Refresher: Drawing random variables

- ▶ Suppose you have a cumulative distribution function  $F_X : \mathcal{X} \subset \mathbb{R} \rightarrow [0, 1]$
- ▶ How can you generate data that is distributed according to  $F$ ?
- ▶ Recall that if  $X \sim F_X$  then  $F_X(x) = \Pr(X \leq x)$
- ▶ But note that  $Z = F_X(X)$  is also a random variable. What is its distribution  $F_Z$ ?
- ▶ Observe that since  $F_X$  is strictly increasing, it has an inverse  $F_X^{-1}$ . So,  
$$F_Z(z) = \Pr(Z \leq z) = \Pr(F_X(X) \leq z) = \Pr(X \leq F_X^{-1}(z)) = F_X(F_X^{-1}(z)) = z$$
- ▶ But this is *exactly* the cdf of the uniform distribution on  $[0,1]$
- ▶  $F_X(X)$  is uniformly distributed!

## Refresher: Drawing random variables

- ▶ Suppose you have a cumulative distribution function  $F_X : \mathcal{X} \subset \mathbb{R} \rightarrow [0, 1]$
- ▶ How can you generate data that is distributed according to  $F$ ?
- ▶ Recall that if  $X \sim F_X$  then  $F_X(x) = \Pr(X \leq x)$
- ▶ But note that  $Z = F_X(X)$  is also a random variable. What is its distribution  $F_Z$ ?
- ▶ Observe that since  $F_X$  is strictly increasing, it has an inverse  $F_X^{-1}$ . So,  
$$F_Z(z) = \Pr(Z \leq z) = \Pr(F_X(X) \leq z) = \Pr(X \leq F_X^{-1}(z)) = F_X(F_X^{-1}(z)) = z$$
- ▶ But this is *exactly* the cdf of the uniform distribution on  $[0,1]$
- ▶  $F_X(X)$  is uniformly distributed!

## Refresher: Drawing random variables

- ▶ Suppose you have a cumulative distribution function  $F_X : \mathcal{X} \subset \mathbb{R} \rightarrow [0, 1]$
- ▶ How can you generate data that is distributed according to  $F$ ?
- ▶ Recall that if  $X \sim F_X$  then  $F_X(x) = \Pr(X \leq x)$
- ▶ But note that  $Z = F_X(X)$  is also a random variable. What is its distribution  $F_Z$ ?
- ▶ Observe that since  $F_X$  is strictly increasing, it has an inverse  $F_X^{-1}$ . So,

$$F_Z(z) = \Pr(Z \leq z) = \Pr(F_X(X) \leq z) = \Pr(X \leq F_X^{-1}(z)) = F_X(F_X^{-1}(z)) = z$$

- ▶ But this is *exactly* the cdf of the uniform distribution on  $[0,1]$
- ▶  $F_X(X)$  is uniformly distributed!

## Refresher: Drawing random variables

- ▶ Suppose you have a cumulative distribution function  $F_X : \mathcal{X} \subset \mathbb{R} \rightarrow [0, 1]$
- ▶ How can you generate data that is distributed according to  $F$ ?
- ▶ Recall that if  $X \sim F_X$  then  $F_X(x) = \Pr(X \leq x)$
- ▶ But note that  $Z = F_X(X)$  is also a random variable. What is its distribution  $F_Z$ ?
- ▶ Observe that since  $F_X$  is strictly increasing, it has an inverse  $F_X^{-1}$ . So,

$$F_Z(z) = \Pr(Z \leq z) = \Pr(F_X(X) \leq z) = \Pr(X \leq F_X^{-1}(z)) = F_X(F_X^{-1}(z)) = z$$

- ▶ But this is *exactly* the cdf of the uniform distribution on  $[0,1]$
- ▶  $F_X(X)$  is uniformly distributed!

## Refresher: Drawing random variables

- ▶ Suppose you have a cumulative distribution function  $F_X : \mathcal{X} \subset \mathbb{R} \rightarrow [0, 1]$
- ▶ How can you generate data that is distributed according to  $F$ ?
- ▶ Recall that if  $X \sim F_X$  then  $F_X(x) = \Pr(X \leq x)$
- ▶ But note that  $Z = F_X(X)$  is also a random variable. What is its distribution  $F_Z$ ?
- ▶ Observe that since  $F_X$  is strictly increasing, it has an inverse  $F_X^{-1}$ . So,

$$F_Z(z) = \Pr(Z \leq z) = \Pr(F_X(X) \leq z) = \Pr(X \leq F_X^{-1}(z)) = F_X(F_X^{-1}(z)) = z$$

- ▶ But this is *exactly* the cdf of the uniform distribution on  $[0,1]$
- ▶  $F_X(X)$  is uniformly distributed!

## Refresher: Drawing random variables

- ▶ Suppose you have a cumulative distribution function  $F_X : \mathcal{X} \subset \mathbb{R} \rightarrow [0, 1]$
- ▶ How can you generate data that is distributed according to  $F$ ?
- ▶ Recall that if  $X \sim F_X$  then  $F_X(x) = \Pr(X \leq x)$
- ▶ But note that  $Z = F_X(X)$  is also a random variable. What is its distribution  $F_Z$ ?
- ▶ Observe that since  $F_X$  is strictly increasing, it has an inverse  $F_X^{-1}$ . So,

$$F_Z(z) = \Pr(Z \leq z) = \Pr(F_X(X) \leq z) = \Pr(X \leq F_X^{-1}(z)) = F_X(F_X^{-1}(z)) = z$$

- ▶ But this is *exactly* the cdf of the uniform distribution on  $[0,1]$
- ▶  $F_X(X)$  is uniformly distributed!

## Refresher: Drawing random variables

- ▶ Suppose you have a cumulative distribution function  $F_X : \mathcal{X} \subset \mathbb{R} \rightarrow [0, 1]$
- ▶ How can you generate data that is distributed according to  $F$ ?
- ▶ Recall that if  $X \sim F_X$  then  $F_X(x) = \Pr(X \leq x)$
- ▶ But note that  $Z = F_X(X)$  is also a random variable. What is its distribution  $F_Z$ ?
- ▶ Observe that since  $F_X$  is strictly increasing, it has an inverse  $F_X^{-1}$ . So,  
$$F_Z(z) = \Pr(Z \leq z) = \Pr(F_X(X) \leq z) = \Pr(X \leq F_X^{-1}(z)) = F_X(F_X^{-1}(z)) = z$$
- ▶ But this is *exactly* the cdf of the uniform distribution on  $[0,1]$
- ▶  $F_X(X)$  is uniformly distributed!

## Refresher: Drawing random variables

- ▶ Suppose you have a cumulative distribution function  $F_X : \mathcal{X} \subset \mathbb{R} \rightarrow [0, 1]$
- ▶ How can you generate data that is distributed according to  $F$ ?
- ▶ Recall that if  $X \sim F_X$  then  $F_X(x) = Pr(X \leq x)$
- ▶ But note that  $Z = F_X(X)$  is also a random variable. What is its distribution  $F_Z$ ?
- ▶ Observe that since  $F_X$  is strictly increasing, it has an inverse  $F_X^{-1}$ . So,  
$$F_Z(z) = Pr(Z \leq z) = Pr(F_X(X) \leq z) = Pr(X \leq F_X^{-1}(z)) = F_X(F_X^{-1}(z)) = z$$
- ▶ But this is *exactly* the cdf of the uniform distribution on  $[0,1]$
- ▶  $F_X(X)$  is uniformly distributed!

## Refresher: Drawing random variables

- ▶ This means that as long as we know how to take pseudo-random draws from a uniform distribution, we can draw random numbers from *any* distribution where we know the cdf
- ▶ Follow the steps:
  1. Draw  $z \sim U[0, 1]$
  2. Find  $x$  such that  $F_X(x) = z$ 

This may require solving a root finding problem if you don't have a formula for  $F_X^{-1}$
  3. Call  $x$  your random draw from  $F$
- ▶ If there are no values of  $X$  that show up with probability 0, then  $F_X$  is strictly increasing and step 2 is always well-defined
- ▶ If there are values of  $X$  that occur with probability 0, then just pick  $x$  as the smallest  $x$  such that  $F_X(x) = z$ .

## Refresher: Drawing random variables

- ▶ This means that as long as we know how to take pseudo-random draws from a uniform distribution, we can draw random numbers from *any* distribution where we know the cdf
- ▶ Follow the steps:
  1. Draw  $z \sim U[0, 1]$
  2. Find  $x$  such that  $F_X(x) = z$ 

This may require solving a root finding problem if you don't have a formula for  $F_X^{-1}$
  3. Call  $x$  your random draw from  $F$
- ▶ If there are no values of  $X$  that show up with probability 0, then  $F_X$  is strictly increasing and step 2 is always well-defined
- ▶ If there are values of  $X$  that occur with probability 0, then just pick  $x$  as the smallest  $x$  such that  $F_X(x) = z$ .

## Refresher: Drawing random variables

- ▶ This means that as long as we know how to take pseudo-random draws from a uniform distribution, we can draw random numbers from *any* distribution where we know the cdf
- ▶ Follow the steps:
  1. Draw  $z \sim U[0, 1]$
  2. Find  $x$  such that  $F_X(x) = z$ 

This may require solving a root finding problem if you don't have a formula for  $F_X^{-1}$
  3. Call  $x$  your random draw from  $F$
- ▶ If there are no values of  $X$  that show up with probability 0, then  $F_X$  is strictly increasing and step 2 is always well-defined
- ▶ If there are values of  $X$  that occur with probability 0, then just pick  $x$  as the smallest  $x$  such that  $F_X(x) = z$ .

# Refresher: Drawing random variables

## Discrete variables

- ▶ This process is particularly straightforward if we have a discrete random variable

$$X \in \{X_1, X_2, \dots, X_n\} \quad \text{where} \quad \Pr(X = X_i) = p_i$$

where  $X_1 < X_2 < \dots < X_n$

- ▶ Observe that we can write the cdf of  $X$  as a cumulative sum:

$$F_i := F(X_i) = \Pr(X \leq X_i) = \sum_{j=1}^i \Pr(X = X_j) = \sum_{j=1}^i p_j$$

- ▶ This means if we draw  $z \sim U[0, 1]$ , “inverting”  $F$  is just a matter of finding the smallest  $i$  such that  $z \leq F_i$
- ▶ Since  $F_i$  are already a sorted vector, this is just searching through a sorted list (fast implementations are easy to find in both Python/Julia)

# Refresher: Drawing random variables

## Discrete variables

- ▶ This process is particularly straightforward if we have a discrete random variable

$$X \in \{X_1, X_2, \dots, X_n\} \quad \text{where} \quad \Pr(X = X_i) = p_i$$

where  $X_1 < X_2 < \dots < X_n$

- ▶ Observe that we can write the cdf of  $X$  as a cumulative sum:

$$F_i := F(X_i) = \Pr(X \leq X_i) = \sum_{j=1}^i \Pr(X = X_j) = \sum_{j=1}^i p_j$$

- ▶ This means if we draw  $z \sim U[0, 1]$ , “inverting”  $F$  is just a matter of finding the smallest  $i$  such that  $z \leq F_i$
- ▶ Since  $F_i$  are already a sorted vector, this is just searching through a sorted list (fast implementations are easy to find in both Python/Julia)

# Refresher: Drawing random variables

## Discrete variables

- ▶ This process is particularly straightforward if we have a discrete random variable

$$X \in \{X_1, X_2, \dots, X_n\} \quad \text{where} \quad \Pr(X = X_i) = p_i$$

where  $X_1 < X_2 < \dots < X_n$

- ▶ Observe that we can write the cdf of  $X$  as a cumulative sum:

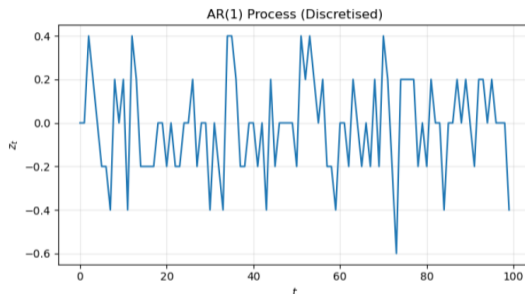
$$F_i := F(X_i) = \Pr(X \leq X_i) = \sum_{j=1}^i \Pr(X = X_j) = \sum_{j=1}^i p_j$$

- ▶ This means if we draw  $z \sim U[0, 1]$ , “inverting”  $F$  is just a matter of finding the smallest  $i$  such that  $z \leq F_i$
- ▶ Since  $F_i$  are already a sorted vector, this is just searching through a sorted list (fast implementations are easy to find in both Python/Julia)

## Refresher: Drawing random variables

```
def simulate_markov(grid, P, T=100, seed=42):  
    np.random.seed(seed)  
    N = len(grid)  
    # Start from middle state  
    i = N // 2  
    X = np.zeros(T)  
    X[0] = grid[i]  
    for t in range(1, T):  
        i = np.random.choice(np.arange(N), p=P[i, :])  
        X[t] = grid[i]  
    return X
```

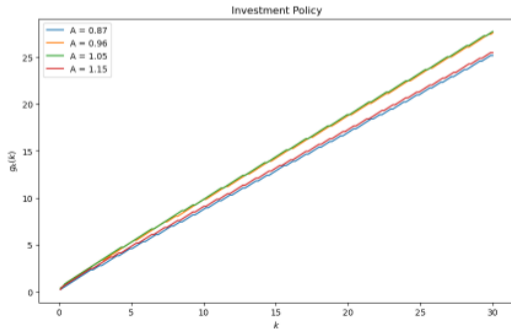
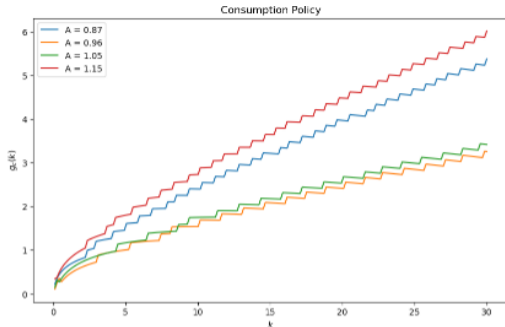
```
grid, P = tauchen(N=7, mu=0.0, rho=0.0, sigma=0.2)  
X = simulate_markov(grid, P, T=100)
```



# Stochastic Neoclassical Growth

## Inner Loop

```
def update_bellman_stochastic(p, m, EV):  
    alpha, beta, delta = p  
    kgrid, Agrid = m["kgrid"], m["Agrid"]  
    Nk, Na = len(kgrid), len(Agrid)  
  
    V = np.zeros((Nk, Na))  
    policy_k = np.zeros((Nk, Na))  
    policy_c = np.zeros((Nk, Na))  
  
    for ai in range(Na):  
        A = Agrid[ai]  
        for ki in range(Nk):  
            k = kgrid[ki]  
            z = A * k**alpha + (1 - delta) * k  
            c = z - kgrid  
            c[c <= 0] = np.nan  
            val = u(c) + beta * EV[:, ai]  
            vmax = np.nanmax(val)  
            pol_index = np.nanargmax(val)  
            V[ki, ai] = vmax  
            policy_k[ki, ai] = kgrid[pol_index]  
            policy_c[ki, ai] = c[pol_index]  
  
    return V, policy_k, policy_c
```



# Stochastic Neoclassical Growth

## Outer Loop

```
def solve_vfi_stochastic(p, m, tol=1e-6, maxiter=1000):  
    alpha, beta, delta = p  
    kgrid, Agrid, V0, P = m["kgrid"], m["Agrid"], m["V0"], m["P"]  
    V = V0.copy()  
    for it in range(maxiter):  
        EV = V @ P  
        Vnew, pol_k, pol_c = update_bellman_stochastic(p, m, EV)  
        err = np.max(np.abs(Vnew - V))  
        if err < tol:  
            print(f"Converged after {it} iterations, err={err:.2e}")  
            break  
        V = Vnew  
    return V, pol_k, pol_c
```

## Section 2

### Estimating Models

## How do we pick our parameters in OLS?

- ▶ When you took econometrics, you learned how to estimate:

$$y_i = X_i\beta + \epsilon_i$$

where  $X_i \in \mathbb{R}^k$ , and  $\beta \in \mathbb{R}^k$

- ▶ Remember that OLS chooses  $\hat{\beta}$  to minimize the sum of squared residuals:

$$\hat{\beta} \in \arg \min_{\beta} \sum_{i=1}^n (y_i - X_i\beta)^2$$

- ▶ It turns out, we can understand OLS as choosing  $\hat{\beta}$  in order to replicate a set of “moments” in our data

A “moment” here means any function mapping data (actual data or predicted data) to a real number. Averages, variances, covariances, etc... are all data moments

- ▶ This called the **method of moments** interpretation of OLS

# How do we pick our parameters in OLS?

- ▶ When you took econometrics, you learned how to estimate:

$$y_i = X_i\beta + \epsilon_i$$

where  $X_i \in \mathbb{R}^k$ , and  $\beta \in \mathbb{R}^k$

- ▶ Remember that OLS chooses  $\hat{\beta}$  to minimize the sum of squared residuals:

$$\hat{\beta} \in \arg \min_{\beta} \sum_{i=1}^n (y_i - X_i\beta)^2$$

- ▶ It turns out, we can understand OLS as choosing  $\hat{\beta}$  in order to replicate a set of “moments” in our data

A “moment” here means any function mapping data (actual data or predicted data) to a real number. Averages, variances, covariances, etc... are all data moments

- ▶ This called the **method of moments** interpretation of OLS

## Method of Moments: OLS

- ▶ To make this easier to see, let's stack up all our dependent observations  $y_i$  in a vector  $y$ , and all our regressors  $X_i$  in a matrix  $X$ .
- ▶ We can write  $SSR(\beta)$  now as a matrix product:

$$SSR(\beta) = (y - X\beta)^T (y - X\beta)$$

- ▶ We can derive the OLS normal equations (first order conditions) by setting  $\frac{\partial SSR}{\partial \beta} = 0$ :

$$X^T (y - X\beta) = 0 \tag{3}$$

- ▶ If we substitute back in  $\hat{\epsilon} = y - X\beta$ , we get the familiar moment conditions:

$$0 = X^T \hat{\epsilon} = \sum_{i=1}^n X_i \hat{\epsilon}_i \iff \text{Cov}(X_i, \hat{\epsilon}_i) = 0 \iff \mathbb{E}[X_i \hat{\epsilon}_i] = 0$$

Note: there are  $k$  different equations here, since  $X_i \in \mathbb{R}^k$

- ▶ OLS chooses  $\beta$  so that the covariance between the predicted residuals and the regressors are zero.
- ▶ This is just a set of  $k$  different "moments"

## Method of Moments: OLS

- ▶ To make this easier to see, let's stack up all our dependent observations  $y_i$  in a vector  $y$ , and all our regressors  $X_i$  in a matrix  $X$ .
- ▶ We can write  $SSR(\beta)$  now as a matrix product:

$$SSR(\beta) = (y - X\beta)^T (y - X\beta)$$

- ▶ We can derive the OLS normal equations (first order conditions) by setting  $\frac{\partial SSR}{\partial \beta} = 0$ :

$$X^T (y - X\beta) = 0 \tag{3}$$

- ▶ If we substitute back in  $\hat{\epsilon} = y - X\beta$ , we get the familiar moment conditions:

$$0 = X^T \hat{\epsilon} = \sum_{i=1}^n X_i \hat{\epsilon}_i \iff \text{Cov}(X_i, \hat{\epsilon}_i) = 0 \iff \mathbb{E}[X_i \hat{\epsilon}_i] = 0$$

Note: there are  $k$  different equations here, since  $X_i \in \mathbb{R}^k$

- ▶ OLS chooses  $\beta$  so that the covariance between the predicted residuals and the regressors are zero.
- ▶ This is just a set of  $k$  different "moments"

## Method of Moments: OLS

- ▶ To make this easier to see, let's stack up all our dependent observations  $y_i$  in a vector  $y$ , and all our regressors  $X_i$  in a matrix  $X$ .
- ▶ We can write  $SSR(\beta)$  now as a matrix product:

$$SSR(\beta) = (y - X\beta)^T (y - X\beta)$$

- ▶ We can derive the OLS normal equations (first order conditions) by setting  $\frac{\partial SSR}{\partial \beta} = 0$ :

$$X^T (y - X\beta) = 0 \tag{3}$$

- ▶ If we substitute back in  $\hat{\epsilon} = y - X\beta$ , we get the familiar moment conditions:

$$0 = X^T \hat{\epsilon} = \sum_{i=1}^n X_i \hat{\epsilon}_i \iff \text{Cov}(X_i, \hat{\epsilon}_i) = 0 \iff \mathbb{E}[X_i \hat{\epsilon}_i] = 0$$

Note: there are  $k$  different equations here, since  $X_i \in \mathbb{R}^k$

- ▶ OLS chooses  $\beta$  so that the covariance between the predicted residuals and the regressors are zero.
- ▶ This is just a set of  $k$  different "moments"

## Method of Moments: OLS

- ▶ To make this easier to see, let's stack up all our dependent observations  $y_i$  in a vector  $y$ , and all our regressors  $X_i$  in a matrix  $X$ .
- ▶ We can write  $SSR(\beta)$  now as a matrix product:

$$SSR(\beta) = (y - X\beta)^T (y - X\beta)$$

- ▶ We can derive the OLS normal equations (first order conditions) by setting  $\frac{\partial SSR}{\partial \beta} = 0$ :

$$X^T (y - X\beta) = 0 \tag{3}$$

- ▶ If we substitute back in  $\hat{\epsilon} = y - X\beta$ , we get the familiar moment conditions:

$$0 = X^T \hat{\epsilon} = \sum_{i=1}^n X_i \hat{\epsilon}_i \iff \text{Cov}(X_i, \hat{\epsilon}_i) = 0 \iff \mathbb{E}[X_i \hat{\epsilon}_i] = 0$$

Note: there are  $k$  different equations here, since  $X_i \in \mathbb{R}^k$

- ▶ OLS chooses  $\beta$  so that the covariance between the predicted residuals and the regressors are zero.
- ▶ This is just a set of  $k$  different "moments"

# Interpreting OLS Method of Moments

- ▶ We can now write our OLS “moment function”  $m : \mathbb{R}^k \rightarrow \mathbb{R}^k$

$$m(\beta) = X^T(y - X\beta)$$

- ▶ Since there is a unique solution to  $m(\beta) = 0$ , we could imagine solving the following problem instead:

$$\hat{\beta} = \arg \min_{\beta} m(\beta)^T m(\beta) \quad (4)$$

- ▶ This is called generalized method of moments: it generalizes OLS to a much broader class of statistical models
  - ▶ In general, you can write IV regressions as a method of moments problem (with a different moment condition)
  - ▶ You can use a moment function which is nonlinear
  - ▶ Crucially, you can use a moment function which involves simulated data from a structural model

# Interpreting OLS Method of Moments

- ▶ We can now write our OLS “moment function”  $m : \mathbb{R}^k \rightarrow \mathbb{R}^k$

$$m(\beta) = X^T(y - X\beta)$$

- ▶ Since there is a unique solution to  $m(\beta) = 0$ , we could imagine solving the following problem instead:

$$\hat{\beta} = \arg \min_{\beta} m(\beta)^T m(\beta) \quad (4)$$

- ▶ This is called generalized method of moments: it generalizes OLS to a much broader class of statistical models
  - ▶ In general, you can write IV regressions as a method of moments problem (with a different moment condition)
  - ▶ You can use a moment function which is nonlinear
  - ▶ Crucially, you can use a moment function which involves simulated data from a structural model

## Stochastic Neoclassical Growth Model: Estimating $\rho$

- ▶ Return to our neoclassical growth model with stochastic productivity

$$\begin{aligned} v(k, A) = \max_{c, k'} & \quad u(c) + \beta \mathbb{E} [v(k', A') | A] \\ \text{s.t.} & \quad c + k' \leq Ak^\alpha + (1 - \delta)k \\ & \quad \log(A') = \rho \log(A) + \epsilon \\ & \quad \epsilon \sim N(0, \sigma) \end{aligned} \tag{2}$$

- ▶ Let's say we want to choose a value of  $\rho$ , but we don't want to just make something up. We want to estimate it to match the data.
- ▶ We know, intuitively, that  $\rho$  controls the persistence of shocks to productivity
  - ▶ If we directly observed  $A_t$ , we could estimate it with a linear regression.
  - ▶ What if, instead, we only observe  $C_t$  and  $Y_t$ ?
  - ▶ Since they both move in step with  $A_t$ , a higher  $\rho$  should give us a more persistent  $C_t$  and  $Y_t$

## Stochastic Neoclassical Growth Model: Estimating $\rho$

- ▶ Return to our neoclassical growth model with stochastic productivity

$$\begin{aligned} v(k, A) = \max_{c, k'} & \quad u(c) + \beta \mathbb{E} [v(k', A') | A] \\ \text{s.t.} & \quad c + k' \leq Ak^\alpha + (1 - \delta)k \\ & \quad \log(A') = \rho \log(A) + \epsilon \\ & \quad \epsilon \sim N(0, \sigma) \end{aligned} \tag{2}$$

- ▶ Let's say we want to choose a value of  $\rho$ , but we don't want to just make something up. We want to estimate it to match the data.
- ▶ We know, intuitively, that  $\rho$  controls the persistence of shocks to productivity
  - ▶ If we directly observed  $A_t$ , we could estimate it with a linear regression.
  - ▶ What if, instead, we only observe  $C_t$  and  $Y_t$ ?
  - ▶ Since they both move in step with  $A_t$ , a higher  $\rho$  should give us a more persistent  $C_t$  and  $Y_t$

# Stochastic Neoclassical Growth Model: Estimating $\rho$

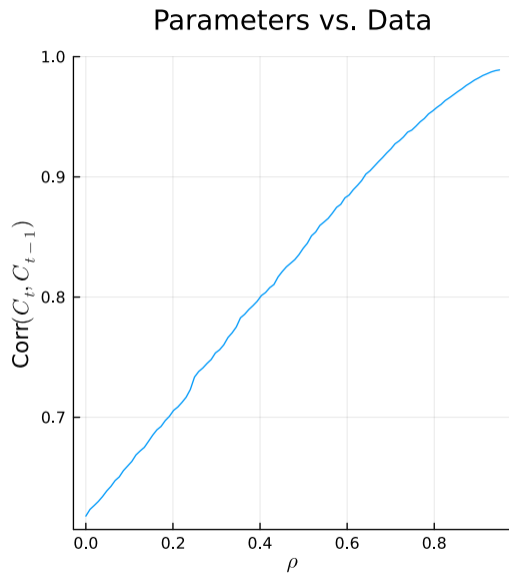
- ▶ Return to our neoclassical growth model with stochastic productivity

$$\begin{aligned} v(k, A) = \max_{c, k'} & \quad u(c) + \beta \mathbb{E} [v(k', A') | A] \\ \text{s.t.} & \quad c + k' \leq Ak^\alpha + (1 - \delta)k \\ & \quad \log(A') = \rho \log(A) + \epsilon \\ & \quad \epsilon \sim N(0, \sigma) \end{aligned} \tag{2}$$

- ▶ Let's say we want to choose a value of  $\rho$ , but we don't want to just make something up. We want to estimate it to match the data.
- ▶ We know, intuitively, that  $\rho$  controls the persistence of shocks to productivity
  - ▶ If we directly observed  $A_t$ , we could estimate it with a linear regression.
  - ▶ What if, instead, we only observe  $C_t$  and  $Y_t$ ?
  - ▶ Since they both move in step with  $A_t$ , a higher  $\rho$  should give us a more persistent  $C_t$  and  $Y_t$

## $\rho$ controls $\text{Corr}(C_t, C_{t-1})$ too

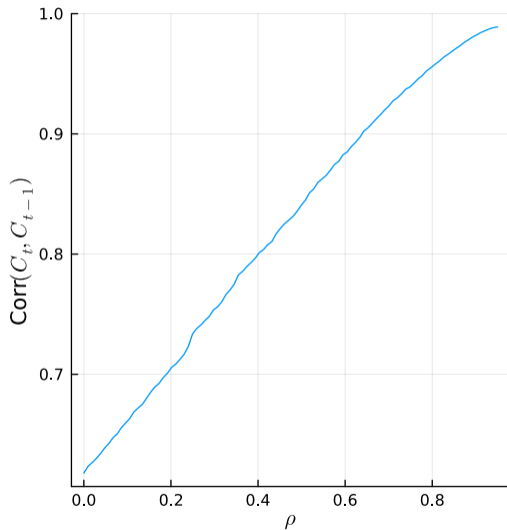
- ▶ What does this tell us?
- ▶ Autocorrelation of consumption increases monotonically with  $\rho$ ,
  - ▶ We should be able to back out a value of  $\rho$  for any “reasonable” value of autocorrelation in consumption that we observe in the data
- ▶ Reasonable means something very specific here:
  - ▶ with the other parameters we've set, we cannot get the model to generate autocorrelation in consumption below around 0.6
  - ▶ It's possible we could if we changed the curvature of the utility function (i.e., made the household less risk averse, so they smooth consumption less)



## $\rho$ controls $\text{Corr}(C_t, C_{t-1})$ too

- ▶ What does this tell us?
- ▶ Autocorrelation of consumption increases monotonically with  $\rho$ ,
  - ▶ We should be able to back out a value of  $\rho$  for any “reasonable” value of autocorrelation in consumption that we observe in the data
- ▶ Reasonable means something very specific here:
  - ▶ with the other parameters we've set, we cannot get the model to generate autocorrelation in consumption below around 0.6
  - ▶ It's possible we could if we changed the curvature of the utility function (i.e., made the household less risk averse, so they smooth consumption less)

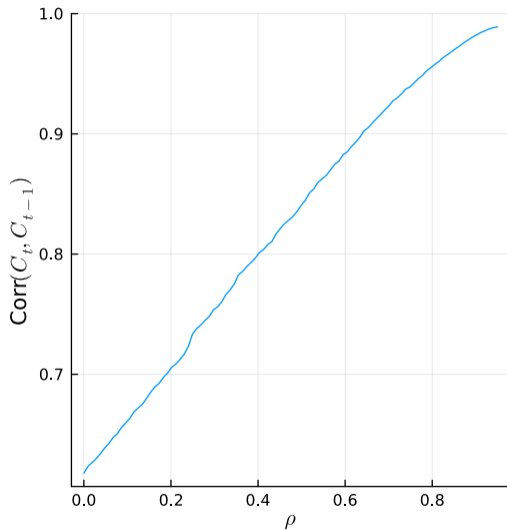
Parameters vs. Data



## $\rho$ controls $\text{Corr}(C_t, C_{t-1})$ too

- ▶ What does this tell us?
- ▶ Autocorrelation of consumption increases monotonically with  $\rho$ ,
  - ▶ We should be able to back out a value of  $\rho$  for any “reasonable” value of autocorrelation in consumption that we observe in the data
- ▶ Reasonable means something very specific here:
  - ▶ with the other parameters we’ve set, we cannot get the model to generate autocorrelation in consumption below around 0.6
  - ▶ It’s possible we could if we changed the curvature of the utility function (i.e, made the household less risk averse, so they smooth consumption less)

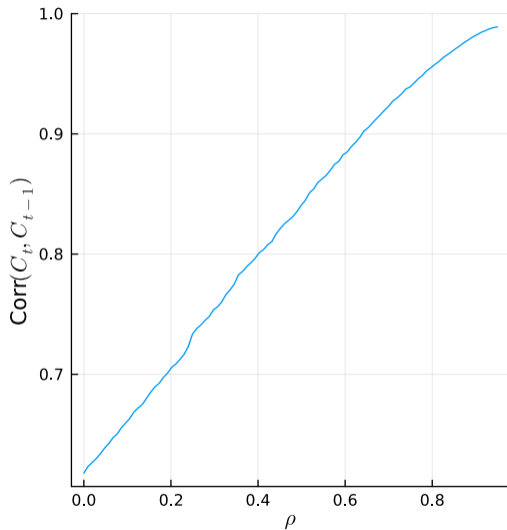
Parameters vs. Data



## $\rho$ controls $\text{Corr}(C_t, C_{t-1})$ too

- ▶ What does this tell us?
- ▶ Autocorrelation of consumption increases monotonically with  $\rho$ ,
  - ▶ We should be able to back out a value of  $\rho$  for any “reasonable” value of autocorrelation in consumption that we observe in the data
- ▶ Reasonable means something very specific here:
  - ▶ with the other parameters we’ve set, we cannot get the model to generate autocorrelation in consumption below around 0.6
  - ▶ It’s possible we could if we changed the curvature of the utility function (i.e, made the household less risk averse, so they smooth consumption less)

Parameters vs. Data



# Estimating $\rho$

- ▶ Suppose you were told (or calculated yourself) that consumption has an autocorrelation coefficient of 0.9

I made this number up – it's just an example

- ▶ How would we choose  $\rho$  so that our model matches the data?
- ▶ Option 1: treat this as a root finding problem
  - ▶ Works fine for just 1 parameters
  - ▶ Runs into issues for many parameters
  - ▶ Simulated moments are either very hard to differentiate, or not differentiable at all, so you will find it hard to calculate a gradient/hessian

# Estimating $\rho$

- ▶ Suppose you were told (or calculated yourself) that consumption has an autocorrelation coefficient of 0.9

I made this number up – it's just an example

- ▶ How would we choose  $\rho$  so that our model matches the data?
- ▶ Option 1: treat this as a root finding problem
  - ▶ Works fine for just 1 parameters
  - ▶ Runs into issues for many parameters
  - ▶ Simulated moments are either very hard to differentiate, or not differentiable at all, so you will find it hard to calculate a gradient/hessian

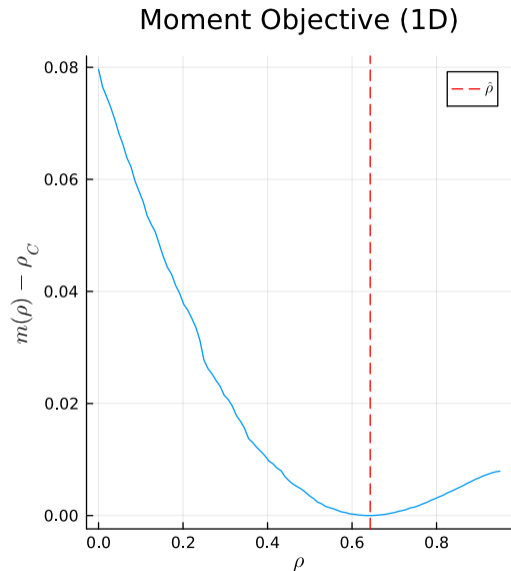
# Estimating $\rho$

- ▶ Option 2: treat this as an optimization problem

$$\hat{\rho} = \arg \min_{\rho} (m(\rho) - \rho_C)^2 \quad (5)$$

where  $m(\rho)$  is the model simulated  $\rho_C$   
(autocorrelation of consumption)

- ▶ This gives the same answer in the 1D case
- ▶ Generalizes better if you want to estimate many parameters at once



# Simulated Method of Moments

- ▶ If we generalize this to several different moments from our data/model, we arrive at the algorithm called **Simulated Method of Moments**
- ▶ Suppose our model has a vector of  $k$  parameters  $\theta \in \mathbb{R}^k$ , which we want to estimate
- ▶ Suppose further, we have a set of  $n \geq k$  moments  $\bar{m} \in \mathbb{R}^n$  (calculated in the data) that are “informative” about the underlying parameters
- ▶ Let  $m : \mathbb{R}^k \rightarrow \mathbb{R}^n$  be the function that
  1. Takes the parameter  $\theta$ , and solves our underlying model
  2. Simulates the model, and
  3. Calculates the model analogue of the moments in the data
- ▶ We choose our estimate  $\hat{\theta}$  to solve

$$\hat{\theta} \in \arg \max_{\theta} (m(\theta) - \bar{m})^T W (m(\theta) - \bar{m}) \quad (6)$$

where  $W$  is a weighting matrix that controls how important each moment is.

For this course, you will generally just set  $W = I$ , and then this is just the sum of squared deviations of model moments from data moments

# Simulated Method of Moments

- ▶ If we generalize this to several different moments from our data/model, we arrive at the algorithm called **Simulated Method of Moments**
- ▶ Suppose our model has a vector of  $k$  parameters  $\theta \in \mathbb{R}^k$ , which we want to estimate
- ▶ Suppose further, we have a set of  $n \geq k$  moments  $\bar{m} \in \mathbb{R}^n$  (calculated in the data) that are “informative” about the underlying parameters
- ▶ Let  $m : \mathbb{R}^k \rightarrow \mathbb{R}^k$  be the function that
  1. Takes the parameter  $\theta$ , and solves our underlying model
  2. Simulates the model, and
  3. Calculates the model analogue of the moments in the data
- ▶ We choose our estimate  $\hat{\theta}$  to solve

$$\hat{\theta} \in \arg \max_{\theta} (m(\theta) - \bar{m})^T W (m(\theta) - \bar{m}) \quad (6)$$

where  $W$  is a weighting matrix that controls how important each moment is.

For this course, you will generally just set  $W = I$ , and then this is just the sum of squared deviations of model moments from data moments

# Simulated Method of Moments

- ▶ If we generalize this to several different moments from our data/model, we arrive at the algorithm called **Simulated Method of Moments**
- ▶ Suppose our model has a vector of  $k$  parameters  $\theta \in \mathbb{R}^k$ , which we want to estimate
- ▶ Suppose further, we have a set of  $n \geq k$  moments  $\bar{m} \in \mathbb{R}^n$  (calculated in the data) that are “informative” about the underlying parameters
- ▶ Let  $m : \mathbb{R}^k \rightarrow \mathbb{R}^k$  be the function that
  1. Takes the parameter  $\theta$ , and solves our underlying model
  2. Simulates the model, and
  3. Calculates the model analogue of the moments in the data

- ▶ We choose our estimate  $\hat{\theta}$  to solve

$$\hat{\theta} \in \arg \max_{\theta} (m(\theta) - \bar{m})^T W (m(\theta) - \bar{m}) \quad (6)$$

where  $W$  is a weighting matrix that controls how important each moment is.

For this course, you will generally just set  $W = I$ , and then this is just the sum of squared deviations of model moments from data moments

## Simulated Method of Moments: Very computationally intensive

- ▶ In general, you have to use a derivative free optimization algorithm (like Nelder-Mead) to solve this problem (if you have even 4 or 5 parameters, this means 100s or 1000s of function evaluations)

- ▶ Even worse, there are local minima everywhere.

Often need to restart optimization from many different starting values and repeat to be sure you've found global min

- ▶ With model solves inside every function evaluation, this could easily run for hours/days/weeks, depending on how large the problem is

Obviously, I won't ask you to solve a problem that has to run for days...

- ▶ You really want to try your best to get blazing fast performance when solving your model

## Simulated Method of Moments: Very computationally intensive

- ▶ In general, you have to use a derivative free optimization algorithm (like Nelder-Mead) to solve this problem (if you have even 4 or 5 parameters, this means 100s or 1000s of function evaluations)

- ▶ Even worse, there are local minima everywhere.

Often need to restart optimization from many different starting values and repeat to be sure you've found global min

- ▶ With model solves inside every function evaluation, this could easily run for hours/days/weeks, depending on how large the problem is

Obviously, I won't ask you to solve a problem that has to run for days...

- ▶ You really want to try your best to get blazing fast performance when solving your model

## Simulated Method of Moments: Very computationally intensive

- ▶ In general, you have to use a derivative free optimization algorithm (like Nelder-Mead) to solve this problem (if you have even 4 or 5 parameters, this means 100s or 1000s of function evaluations)

- ▶ Even worse, there are local minima everywhere.

Often need to restart optimization from many different starting values and repeat to be sure you've found global min

- ▶ With model solves inside every function evaluation, this could easily run for hours/days/weeks, depending on how large the problem is

Obviously, I won't ask you to solve a problem that has to run for days...

- ▶ You really want to try your best to get blazing fast performance when solving your model

# Simulated Method of Moments: Very powerful

- ▶ All you need are “informative” moments
  - ▶ If  $f(\theta)$  is your objective, then you need  $Df(\theta^*)$  to have full rank (be invertible) at the true parameters
- ▶ Any statistic can be “informative” if it captures the right thing (although choosing moments is a bit of an art)

You can use aggregate statistics, which is great for working with confidential microdata

- ▶ In principle you can use regression coefficients that are badly identified, so long as the identification problem you're worried about in the data is also present in the model

Example: labor search models with unobserved heterogeneity in human capital. There is selection on who is hired out of unemployment (since employers see your skill)

As long as you model the unobserved heterogeneity in human capital, you can make use of a regression of earnings just out of unemployment on duration of unemployment

# Simulated Method of Moments: Very powerful

- ▶ All you need are “informative” moments
  - ▶ If  $f(\theta)$  is your objective, then you need  $Df(\theta^*)$  to have full rank (be invertible) at the true parameters
- ▶ Any statistic can be “informative” if it captures the right thing (although choosing moments is a bit of an art)

You can use aggregate statistics, which is great for working with confidential microdata

- ▶ In principle you can use regression coefficients that are badly identified, so long as the identification problem you're worried about in the data is also present in the model

Example: labor search models with unobserved heterogeneity in human capital. There is selection on who is hired out of unemployment (since employers see your skill)

As long as you model the unobserved heterogeneity in human capital, you can make use of a regression of earnings just out of unemployment on duration of unemployment

# Simulated Method of Moments: Very powerful

- ▶ All you need are “informative” moments
  - ▶ If  $f(\theta)$  is your objective, then you need  $Df(\theta^*)$  to have full rank (be invertible) at the true parameters
- ▶ Any statistic can be “informative” if it captures the right thing (although choosing moments is a bit of an art)

You can use aggregate statistics, which is great for working with confidential microdata

- ▶ In principle you can use regression coefficients that are badly identified, so long as the identification problem you're worried about in the data is also present in the model

Example: labor search models with unobserved heterogeneity in human capital. There is selection on who is hired out of unemployment (since employers see your skill)

As long as you model the unobserved heterogeneity in human capital, you can make use of a regression of earnings just out of unemployment on duration of unemployment

# Summary

- ▶ Saw today that you can simulate your model fairly easily once you've solved it
  - Just need to repeatedly apply the policy functions you've solved for
- ▶ Simulating the model can be extremely helpful for estimating your model to match the data
  - ▶ Simulated method of moments is a very powerful tool for estimating structural models
  - ▶ Only requires that parameters be “informative” about something in the data to estimate them
  - ▶ It's generally quite tricky to pick sensible moments for estimation (something of an art)
- ▶ Next week, we'll return to see how we can use these models to run policy counterfactuals