

# WHAT DO DATA ON MILLIONS OF U.S. WORKERS REVEAL ABOUT LIFE CYCLE INCOME RISK?<sup>1</sup>

**Fatih Guvenen**

Minnesota and NBER

**Fatih Karahan**

New York Fed

**Serdar Ozkan**

Toronto

**Jae Song**

SSA

**March 26, 2015**

---

<sup>1</sup>The findings and conclusions expressed are solely those of the authors and do not represent the views of Federal Reserve Board, Federal Reserve Bank of New York or SSA.

# THREE QUESTIONS

1. What does the distribution of earnings shocks look like?

# THREE QUESTIONS

1. What does the distribution of earnings shocks look like?
  - Is it approximately **lognormal**?
  - Symmetric or skewed? any excess kurtosis?

# THREE QUESTIONS

1. What does the distribution of earnings shocks look like?
  - Is it approximately **lognormal**?
  - Symmetric or skewed? any excess kurtosis?
2. How do these properties vary:
  - A. over the **life cycle**?
  - B. across **income groups**?

# THREE QUESTIONS

1. What does the distribution of earnings shocks look like?
  - Is it approximately **lognormal**?
  - Symmetric or skewed? any excess kurtosis?
2. How do these properties vary:
  - A. over the **life cycle**?
  - B. across **income groups**?
3. Dynamics of earnings:
  - A. Do positive and negative shocks have similar persistence?
  - B. Do large and small shocks have similar persistence?

# THIS PAPER

## Existing work:

1. Small survey-based data sets, e.g., the PSID
  - between 500 to 2000 individuals per year

## This paper:

# THIS PAPER

## Existing work:

1. Small survey-based data sets, e.g., the PSID
  - between 500 to 2000 individuals per year

## This paper:

1. Large and clean administrative data set
  - as many as 5,000,000 individuals per year.

# THIS PAPER

## Existing work:

1. Small survey-based data sets, e.g., the PSID
  - between 500 to 2000 individuals per year
2. Employ covariance matrix estimation (CME), developed for a data-constrained environment
  - *Notable exceptions:* Meghir and Pistaferri (2004), Browning et al (2010), and Altonji et al (2013)

## This paper:

1. Large and clean administrative data set
  - as many as 5,000,000 individuals per year.



# THIS PAPER

## Existing work:

1. Small survey-based data sets, e.g., the PSID
  - between 500 to 2000 individuals per year
2. Employ covariance matrix estimation (CME), developed for a data-constrained environment
  - *Notable exceptions:* Meghir and Pistaferri (2004), Browning et al (2010), and Altonji et al (2013)

## This paper:

1. Large and clean administrative data set
  - as many as 5,000,000 individuals per year.
2. Move beyond CME and target **economically significant** moments.

# COVARIANCE MATRIX ESTIMATION: A RECAP

1. Specify a parametric income process, e.g.:

$$y_t^i = \alpha^i + z_t^i + \varepsilon_t^i$$

$$z_t^i = \rho z_{t-1}^i + \eta_t^i$$

# COVARIANCE MATRIX ESTIMATION: A RECAP

1. Specify a parametric income process, e.g.:

$$y_t^i = \alpha^i + z_t^i + \varepsilon_t^i$$

$$z_t^i = \rho z_{t-1}^i + \eta_t^i$$

2. Derive the theoretical autocovariances of income implied by this specification:

$$\text{var}_i(y_t^i) = \sigma_\alpha^2 + \text{var}_i(z_t^i) + \sigma_\varepsilon^2,$$

$$\text{var}_i(z_t^i) = \sum_{s=1}^t \rho^{2s} \sigma_\eta^2,$$

$$\text{cov}(y_t^i, y_{t+n}^i) = \sigma_\alpha^2 + \rho^n \text{var}_i(z_t^i).$$

# COVARIANCE MATRIX ESTIMATION

3. Construct the empirical counterpart of the covariance matrix:

$$\begin{bmatrix} \text{var}_i(y_1^i) & & & & & \\ & \vdots & & & & \\ & & \text{var}_i(y_2^i) & & & \\ & & & \dots & & \\ & & & & \ddots & \\ \text{cov}(y_1^i, y_t^i) & & & \dots & & \text{var}_i(y_t^i) \\ & \vdots & & & & & \\ \text{cov}(y_1^i, y_T^i) & & & & & \dots & \text{var}_i(y_T^i) \end{bmatrix}$$

## COVARIANCE MATRIX ESTIMATION

- Construct the empirical counterpart of the covariance matrix:

$$\begin{bmatrix} \text{var}_i(y_1^i) & & & & \\ \vdots & \text{var}_i(y_2^i) & & & \\ \vdots & \dots & \ddots & & \\ \text{cov}(y_1^i, y_t^i) & \dots & & \text{var}_i(y_t^i) & \\ \vdots & & & & \ddots \\ \text{cov}(y_1^i, y_T^i) & & & & \text{var}_i(y_T^i) \end{bmatrix}$$

- Choose  $(\rho, \sigma_\alpha^2, \sigma_\eta^2, \sigma_\varepsilon^2)$  to bring the theoretical covariance matrix as close to its empirical part as possible.

# MOVING BEYOND THE COVARIANCE MATRIX

- Covariances lump many disparate features of the data into one statistic.
  - Discards lots of useful information.

# MOVING BEYOND THE COVARIANCE MATRIX

- Covariances lump many disparate features of the data into one statistic.
  - Discards lots of useful information.
- Also ignores **higher-order moments**, which we find to be very important.

# MOVING BEYOND THE COVARIANCE MATRIX

- Covariances lump many disparate features of the data into one statistic.
  - Discards lots of useful information.
- Also ignores **higher-order moments**, which we find to be very important.
- With CME, selecting among rejected models is very hard:
  - moments that are missed **do not have clear economic interpretations**.



# THIS PAPER

Uses a unique, confidential, and very large administrative dataset to:

# THIS PAPER

Uses a unique, confidential, and very large administrative dataset to:

1. **Document** new empirical facts on life-cycle earnings dynamics

# THIS PAPER

Uses a unique, confidential, and very large administrative dataset to:

1. **Document** new empirical facts on life-cycle earnings dynamics
2. **Estimate** lifecycle earnings dynamics
  - by matching **economically important moments** (as opposed to the “covariance matrix of income residuals”)

# THIS PAPER

Uses a unique, confidential, and very large administrative dataset to:

1. **Document** new empirical facts on life-cycle earnings dynamics
2. **Estimate** lifecycle earnings dynamics
  - by matching **economically important moments** (as opposed to the “covariance matrix of income residuals”)
3. Provide a reliable “**user’s guide**” for earnings process specifications.

# NEW EMPIRICAL FACTS

# DATA: SSA MASTER EARNINGS FILE

- **Representative sample of US males** covering 34 years:  
1978 to 2011

## DATA: SSA MASTER EARNINGS FILE

- Representative sample of US males covering 34 years: 1978 to 2011
- Salary and wage workers (from W-2 forms)

## DATA: SSA MASTER EARNINGS FILE

- Representative sample of US males covering 34 years: 1978 to 2011
- Salary and wage workers (from W-2 forms)
- Individuals aged 25–60



## DATA: SSA MASTER EARNINGS FILE

- Representative sample of US males covering 34 years: 1978 to 2011
- Salary and wage workers (from W-2 forms)
- Individuals aged 25–60
- Key Advantages:
  - Very large sample size (200+ million observations)
  - No survey response error
  - No sample attrition
  - No top-coding

# FOUR SETS OF EMPIRICAL FACTS

1. Average income growth over the life cycle

## FOUR SETS OF EMPIRICAL FACTS

1. Average income growth over the life cycle
2. Cross-sectional moments of earnings growth

## FOUR SETS OF EMPIRICAL FACTS

1. Average income growth over the life cycle
2. Cross-sectional moments of earnings growth
3. Short- and long-run dynamics of income growth

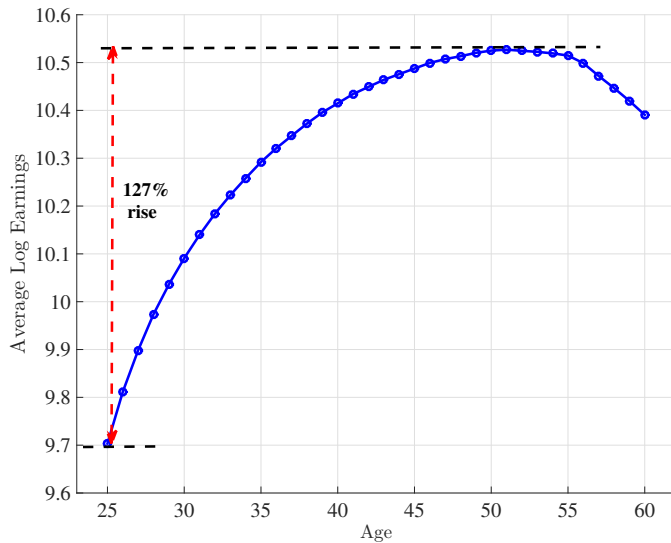
## FOUR SETS OF EMPIRICAL FACTS

1. Average income growth over the life cycle
2. Cross-sectional moments of earnings growth
3. Short- and long-run dynamics of income growth
4. Scarring Effects of Long-Term Unemployment
5. Distribution of Lifetime Income (skip today)

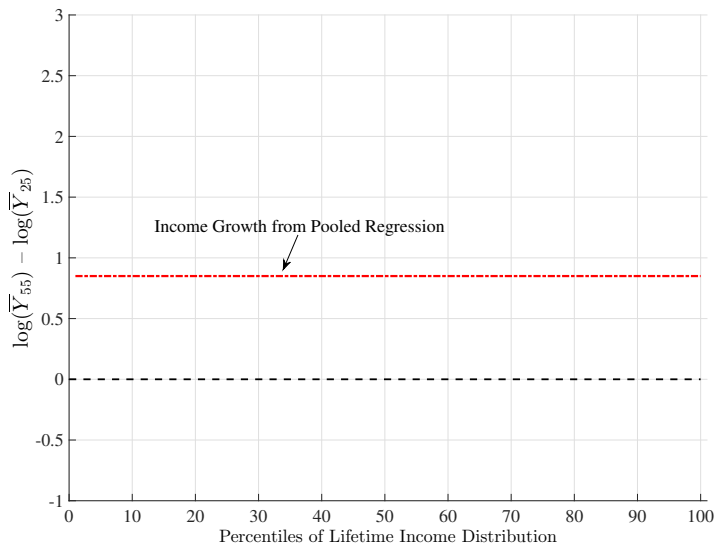
# FOUR SETS OF EMPIRICAL FACTS

1. Average Income growth over the life cycle
- 2.
- 3.
- 4.

# I. AGE PROFILE OF LABOR INCOME

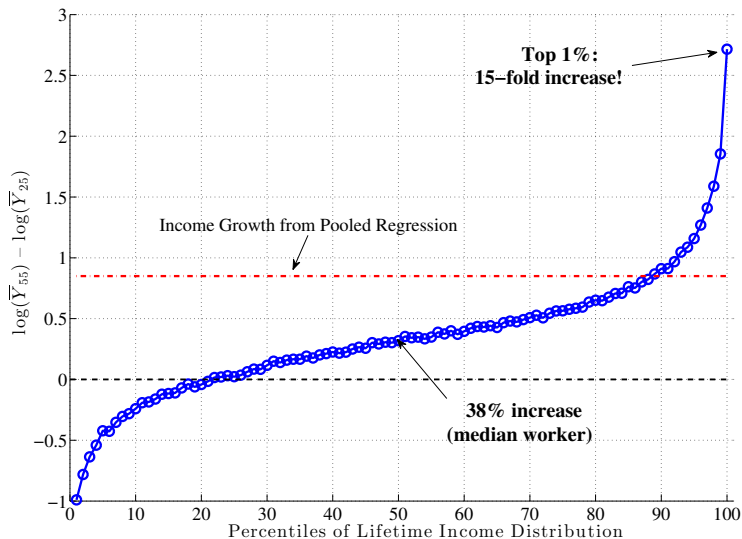


# I. INCOME GROWTH OVER LIFE CYCLE

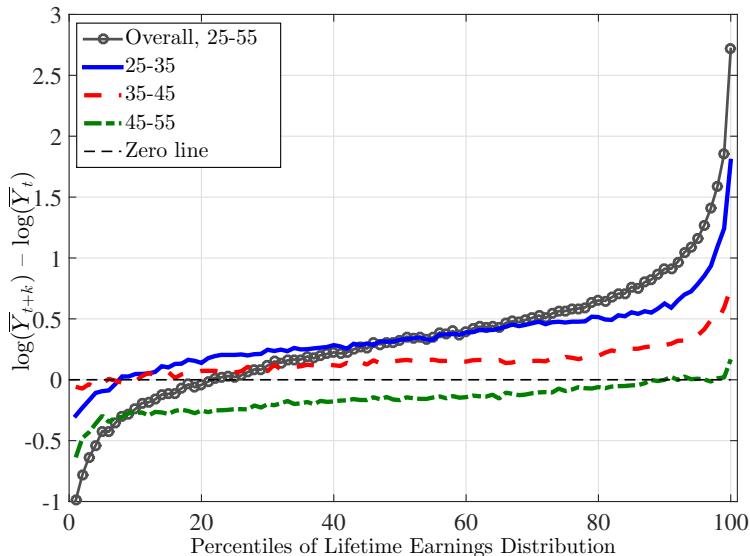




# I. INCOME GROWTH OVER LIFE CYCLE



# I. INCOME GROWTH OVER LIFE CYCLE

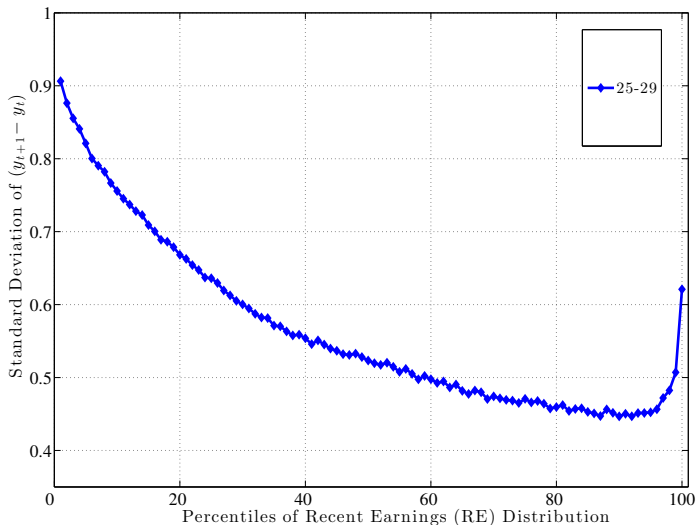


# FOUR SETS OF EMPIRICAL FACTS

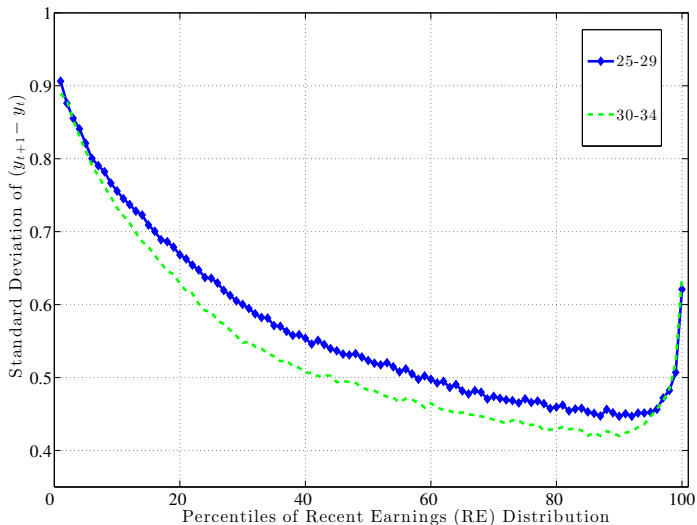
1. Income growth over the life cycle
2. Cross-sectional moments of earnings growth:  $y_{t+k} - y_t$
- 3.
- 4.

# Standard Deviation

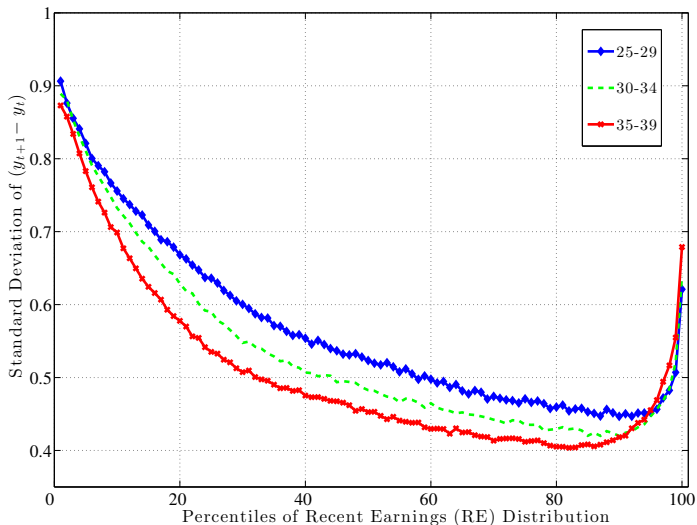
## II.A STANDARD DEVIATION OF $y_{t+1} - y_t$



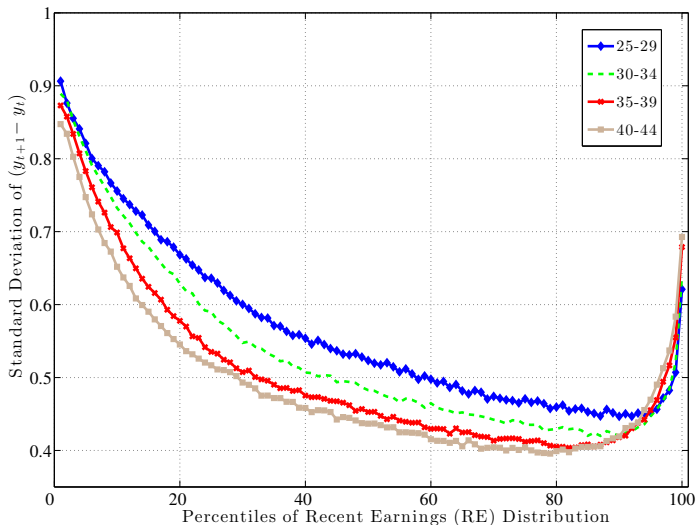
## II.A STANDARD DEVIATION OF $y_{t+1} - y_t$



## II.A STANDARD DEVIATION OF $y_{t+1} - y_t$

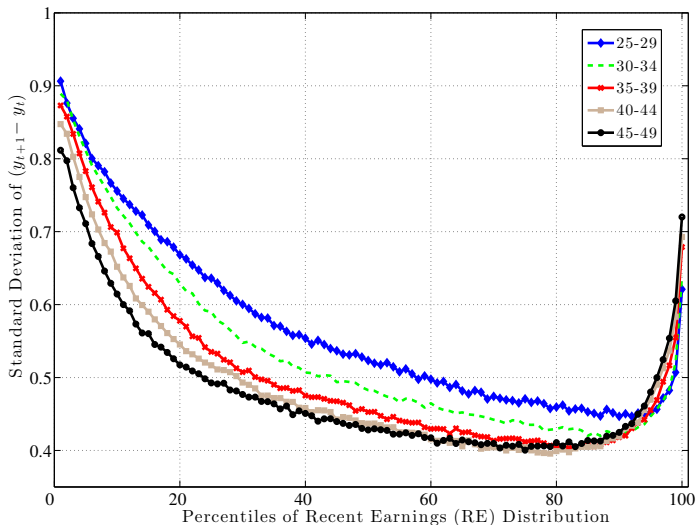


## II.A STANDARD DEVIATION OF $y_{t+1} - y_t$

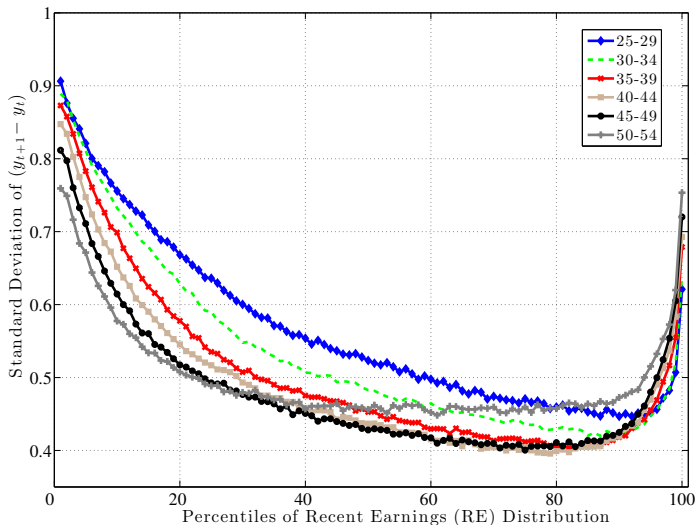




## II.A STANDARD DEVIATION OF $y_{t+1} - y_t$

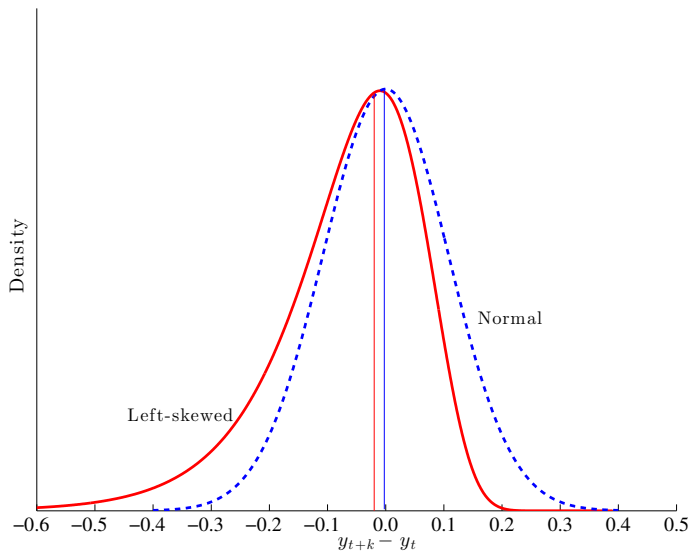


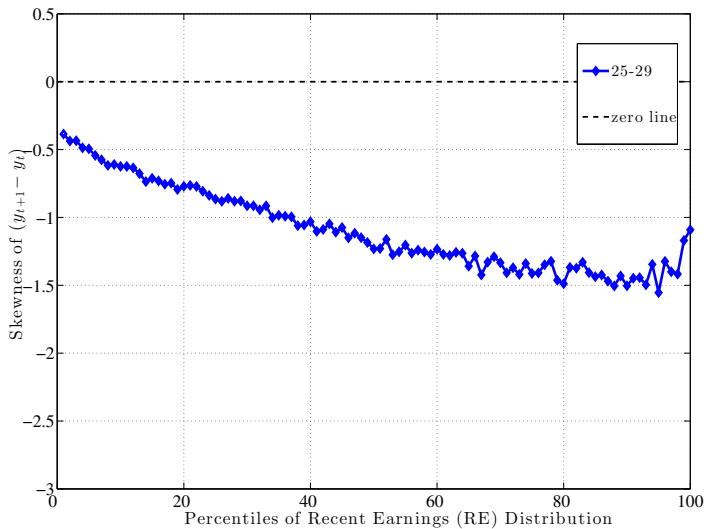
## II.A STANDARD DEVIATION OF $y_{t+1} - y_t$

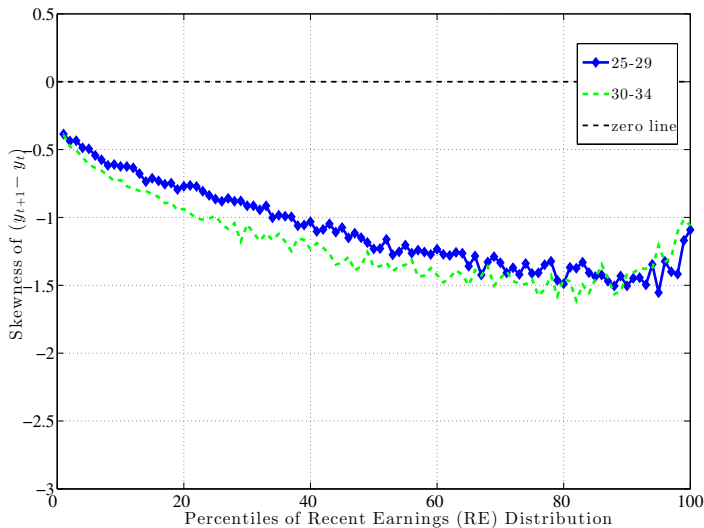


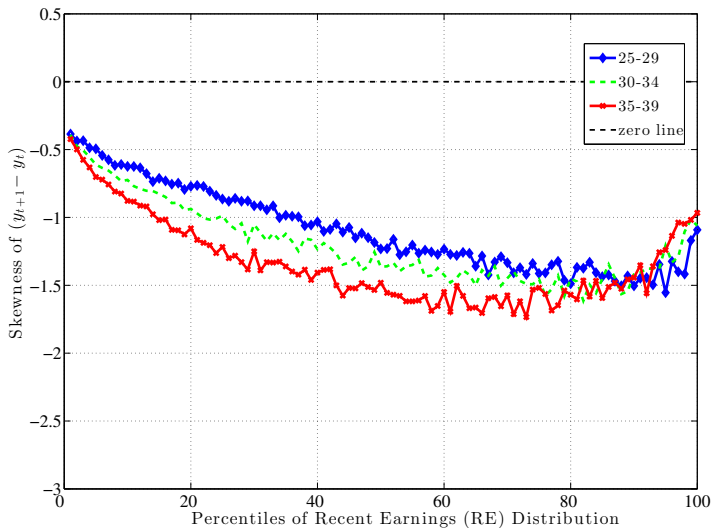
Skewness

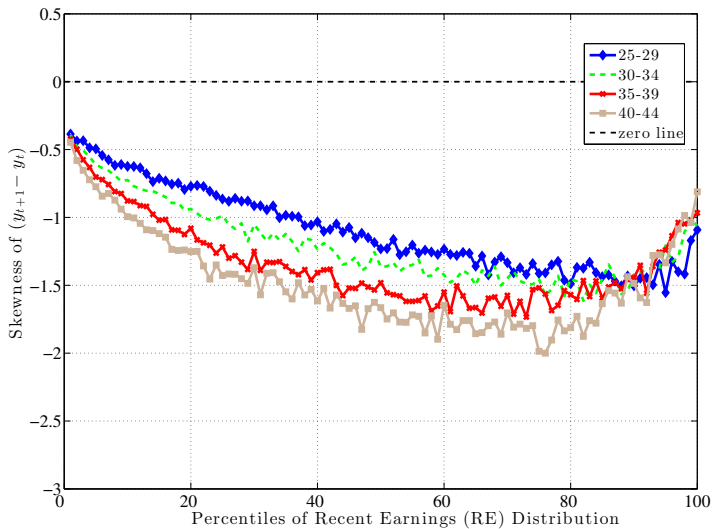
# LEFT-SKEWNESS



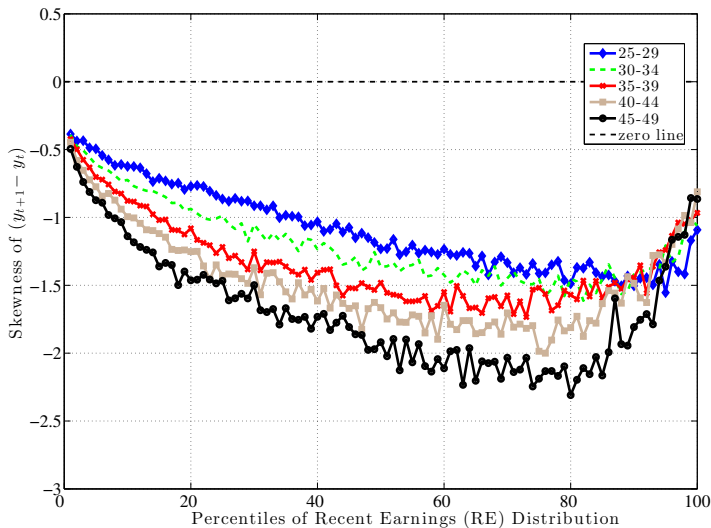
II.B SKEWNESS OF  $y_{t+1} - y_t$ 

II.B SKEWNESS OF  $y_{t+1} - y_t$ 

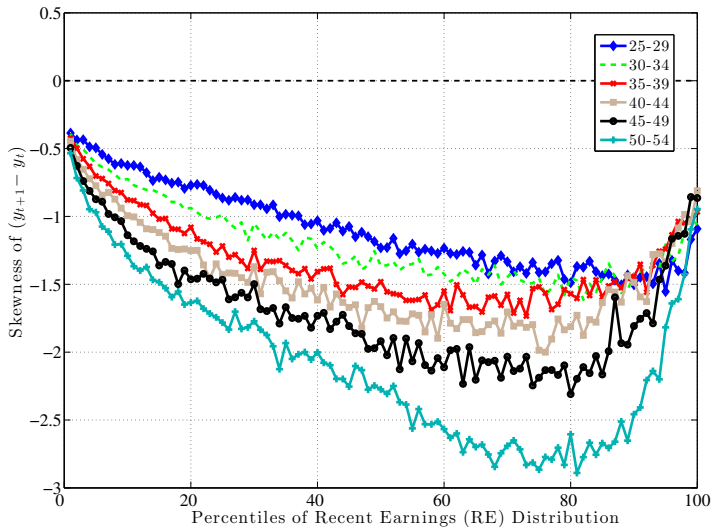
II.B SKEWNESS OF  $y_{t+1} - y_t$ 

II.B SKEWNESS OF  $y_{t+1} - y_t$ 



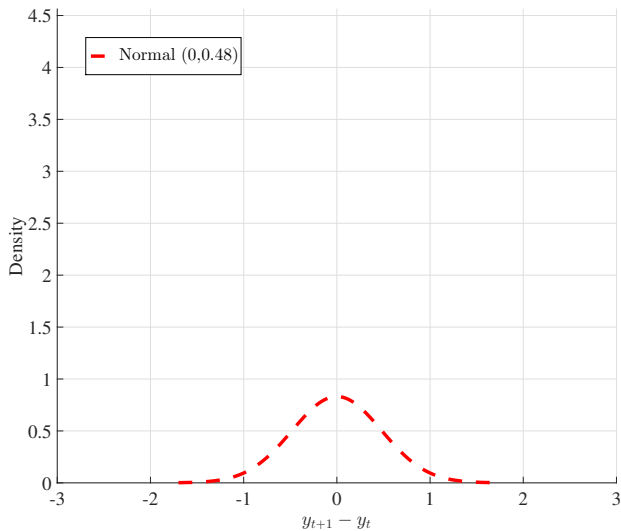
II.B SKEWNESS OF  $y_{t+1} - y_t$ 

## II.B SKEWNESS OF $y_{t+1} - y_t$

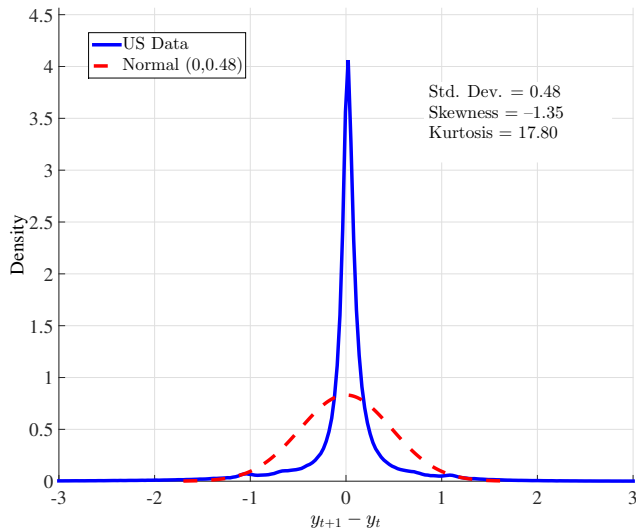


Kurtosis

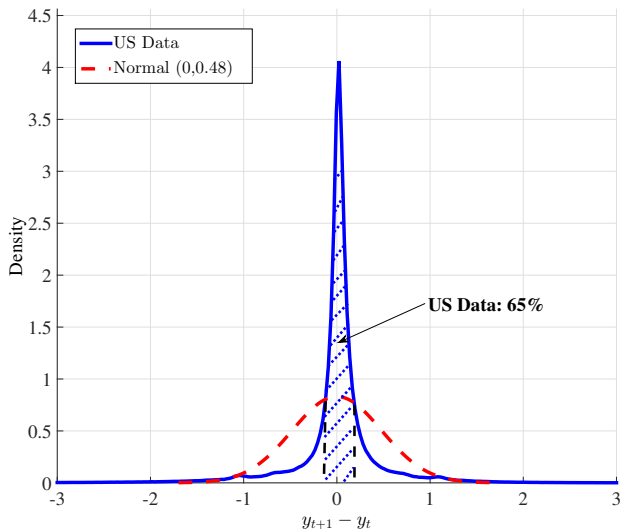
## II.C HISTOGRAM OF $y_{t+1} - y_t$



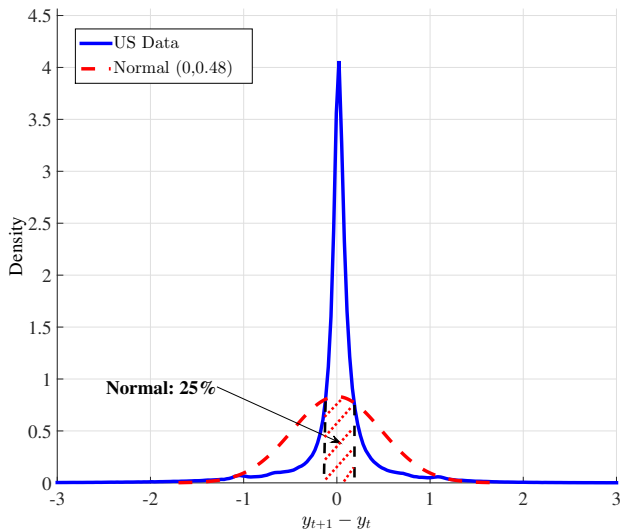
## II.C HISTOGRAM OF $y_{t+1} - y_t$



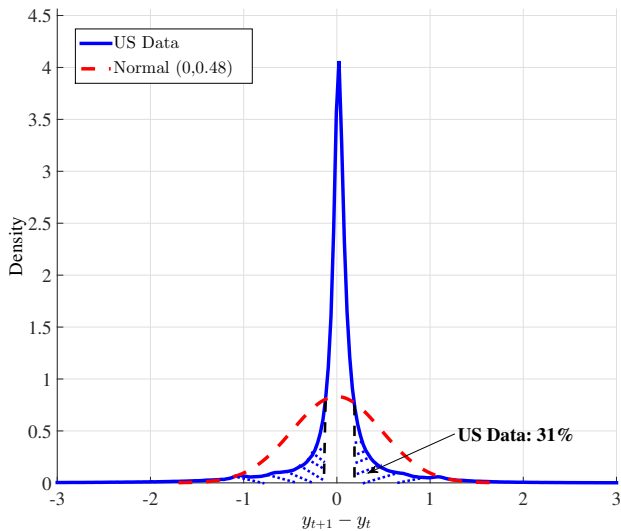
## II.C CENTER OF HISTOGRAM: $[-0.12, 0.19]$



## II.C CENTER OF HISTOGRAM: $[-0.12, 0.19]$

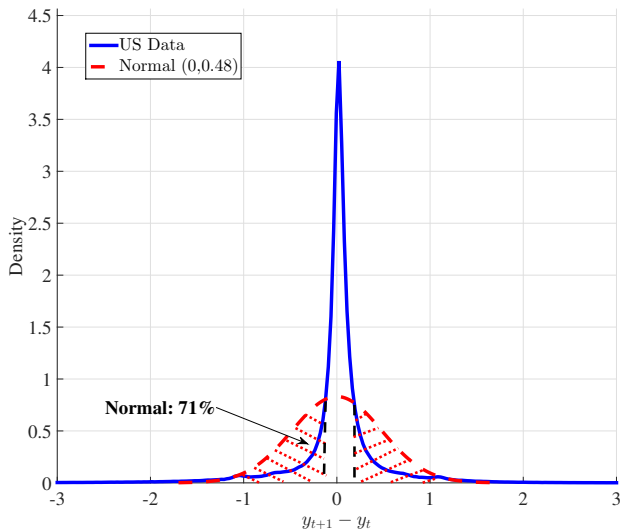


## II.C SHOULDERS OF HISTOGRAM

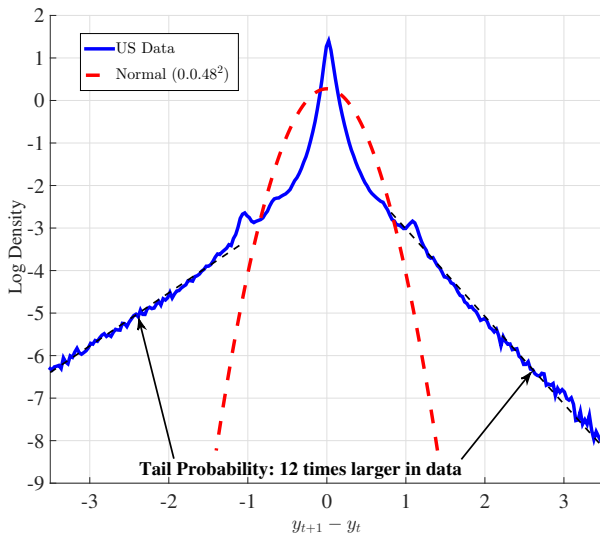




## II.C SHOULDERS OF HISTOGRAM



# I. TAILS OF HISTOGRAM

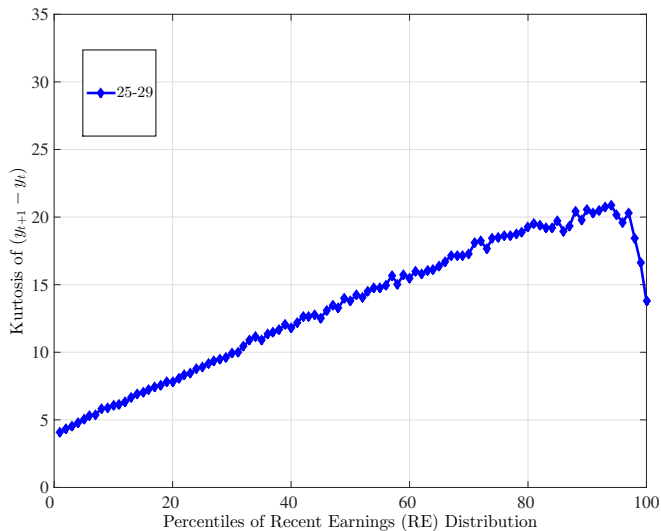


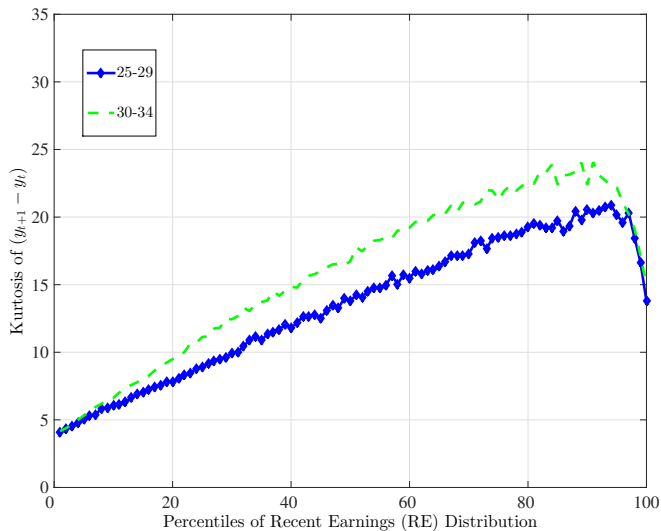
## II.C DISTRIBUTION OF INCOME CHANGES

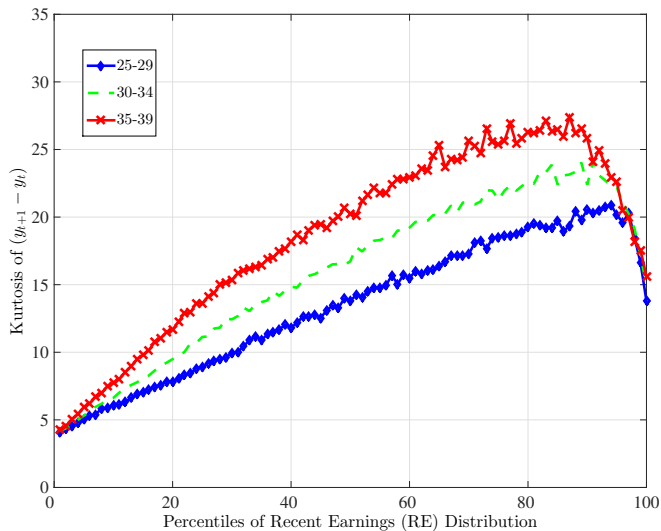
$x :$	$\text{Prob}( y_{t+1}^i - y_t^i  < x)$		
	Data*	$\mathcal{N}(0, 0.48^2)$	Ratio
<b>0.05</b>	<b>0.35</b>	<b>0.08</b>	<b>4.38</b>
0.10	0.54	0.16	3.38
0.20	0.71	0.32	2.23
0.50	0.86	0.70	1.22
1.00	0.94	0.96	0.98

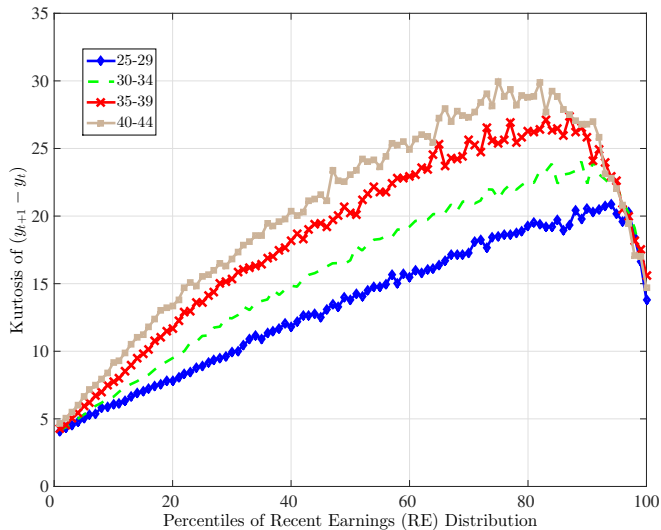
## II.C DISTRIBUTION OF INCOME CHANGES

$x :$	$\text{Prob}( y_{t+1}^i - y_t^i  < x)$		
	Data*	$\mathcal{N}(0, 0.48^2)$	Ratio
0.05	0.35	0.08	4.38
<b>0.10</b>	<b>0.54</b>	<b>0.16</b>	<b>3.38</b>
0.20	0.71	0.32	2.23
0.50	0.86	0.70	1.22
1.00	0.94	0.96	0.98

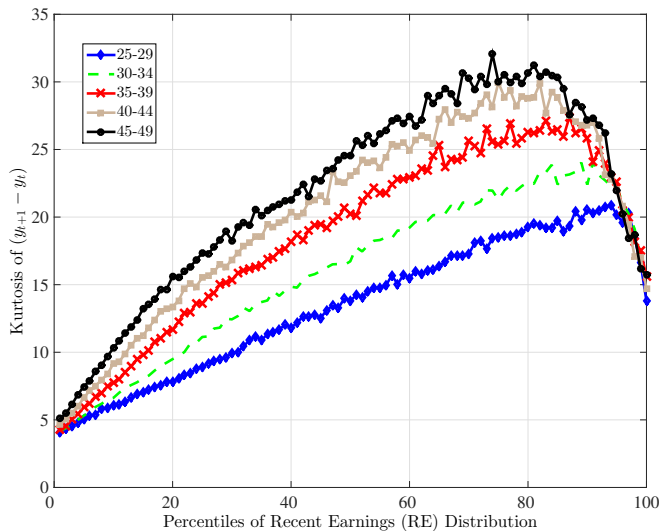
II.C KURTOSIS OF  $y_{t+1} - y_t$ 

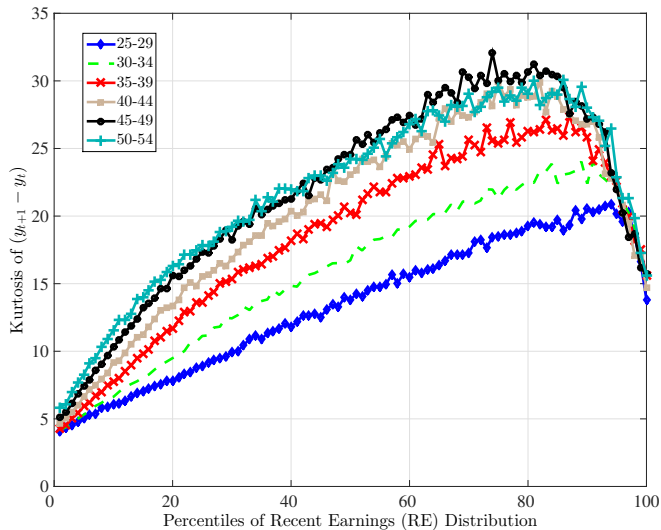
II.C KURTOSIS OF  $y_{t+1} - y_t$ 

II.C KURTOSIS OF  $y_{t+1} - y_t$ 

II.C KURTOSIS OF  $y_{t+1} - y_t$ 



II.C KURTOSIS OF  $y_{t+1} - y_t$ 

II.C KURTOSIS OF  $y_{t+1} - y_t$ 

## RISK PREMIUM: SKEWNESS AND KURTOSIS

Let  $\tilde{\delta}$  be a static gamble. And  $\pi$  is the risk premium to avoid it:

$$U(c \times (1 - \pi)) = \mathbb{E} \left[ U(c \times (1 + \tilde{\delta})) \right].$$

## RISK PREMIUM: SKEWNESS AND KURTOSIS

Let  $\tilde{\delta}$  be a static gamble. And  $\pi$  is the risk premium to avoid it:

$$U(c \times (1 - \pi)) = \mathbb{E} \left[ U(c \times (1 + \tilde{\delta})) \right].$$

<i>Gamble:</i>	Risk Premium ( $\pi$ )	
	$\tilde{\delta}^A$	$\tilde{\delta}^B$
Mean	0.0	
Standard Deviation	0.10	
Skewness	0.0	
Excess Kurtosis	0.0	
Premium	.	

## RISK PREMIUM: SKEWNESS AND KURTOSIS

Let  $\tilde{\delta}$  be a static gamble. And  $\pi$  is the risk premium to avoid it:

$$U(c \times (1 - \pi)) = \mathbb{E} \left[ U(c \times (1 + \tilde{\delta})) \right].$$

<i>Gamble:</i>	Risk Premium ( $\pi$ )	
	$\tilde{\delta}^A$	$\tilde{\delta}^B$
Mean	0.0	
Standard Deviation	0.10	
Skewness	0.0	
Excess Kurtosis	0.0	
Premium	<b>4.88%</b>	

## RISK PREMIUM: SKEWNESS AND KURTOSIS

Let  $\tilde{\delta}$  be a static gamble. And  $\pi$  is the risk premium to avoid it:

$$U(c \times (1 - \pi)) = \mathbb{E} \left[ U(c \times (1 + \tilde{\delta})) \right].$$

<i>Gamble:</i>	Risk Premium ( $\pi$ )	
	$\tilde{\delta}^A$	$\tilde{\delta}^B$
Mean	0.0	0.0
Standard Deviation	0.10	0.10
Skewness	0.0	.
Excess Kurtosis	0.0	.
Premium	<b>4.88%</b>	.

## RISK PREMIUM: SKEWNESS AND KURTOSIS

Let  $\tilde{\delta}$  be a static gamble. And  $\pi$  is the risk premium to avoid it:

$$U(c \times (1 - \pi)) = \mathbb{E} \left[ U(c \times (1 + \tilde{\delta})) \right].$$

<i>Gamble:</i>	Risk Premium ( $\pi$ )	
	$\tilde{\delta}^A$	$\tilde{\delta}^B$
Mean	0.0	0.0
Standard Deviation	0.10	0.10
Skewness	0.0	-2.0
Excess Kurtosis	0.0	27.0
Premium	4.88%	.

## RISK PREMIUM: SKEWNESS AND KURTOSIS

Let  $\tilde{\delta}$  be a static gamble. And  $\pi$  is the risk premium to avoid it:

$$U(c \times (1 - \pi)) = \mathbb{E} \left[ U(c \times (1 + \tilde{\delta})) \right].$$

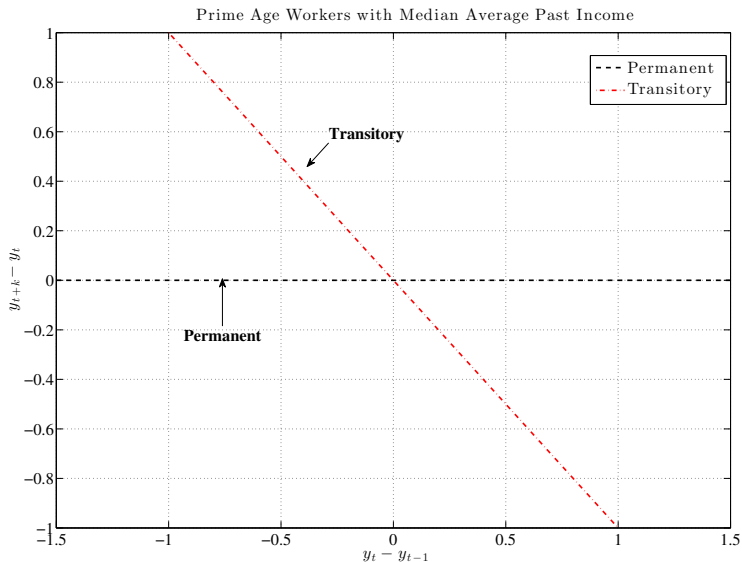
<i>Gamble:</i>	Risk Premium ( $\pi$ )	
	$\tilde{\delta}^A$	$\tilde{\delta}^B$
Mean	0.0	0.0
Standard Deviation	0.10	0.10
Skewness	0.0	-2.0
Excess Kurtosis	0.0	27.0
Premium	4.88%	<b>22.15%</b>



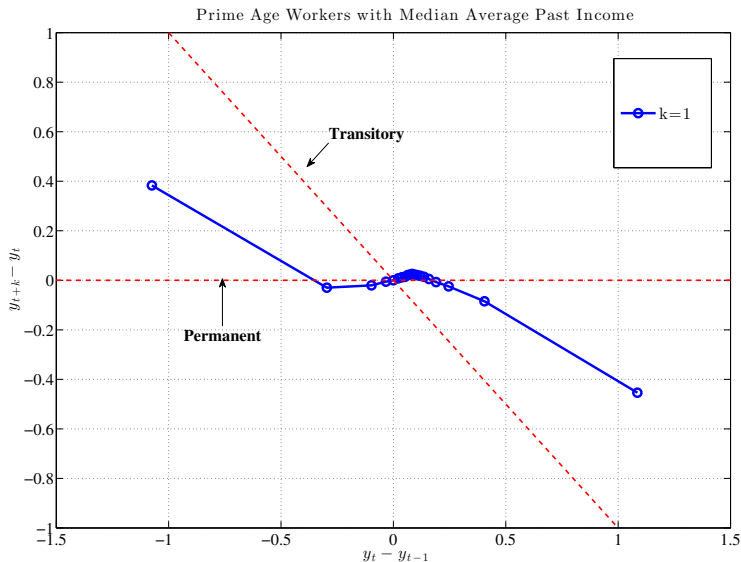
# FOUR SETS OF EMPIRICAL FACTS

1. Average income growth over the life cycle
2. Cross-sectional moments of earnings growth:  $y_{t+k} - y_t$
3. Short- and long-run dynamics of income growth
- 4.

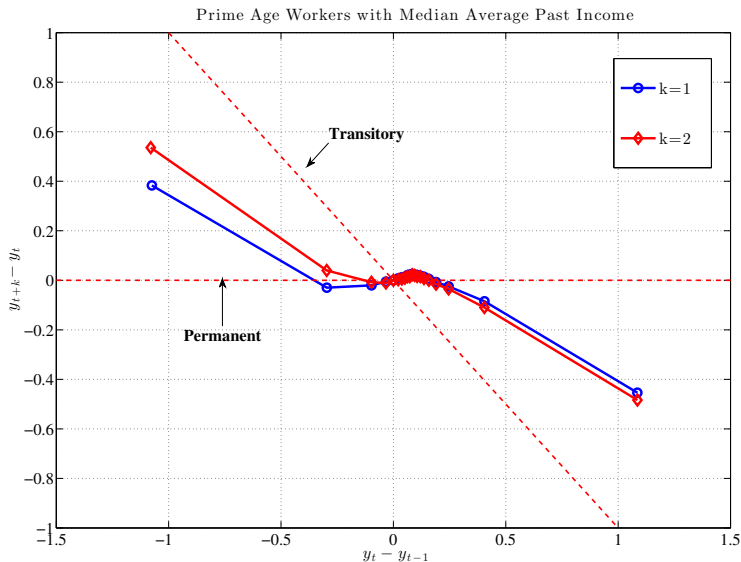
# IMPULSE RESPONSE FUNCTIONS



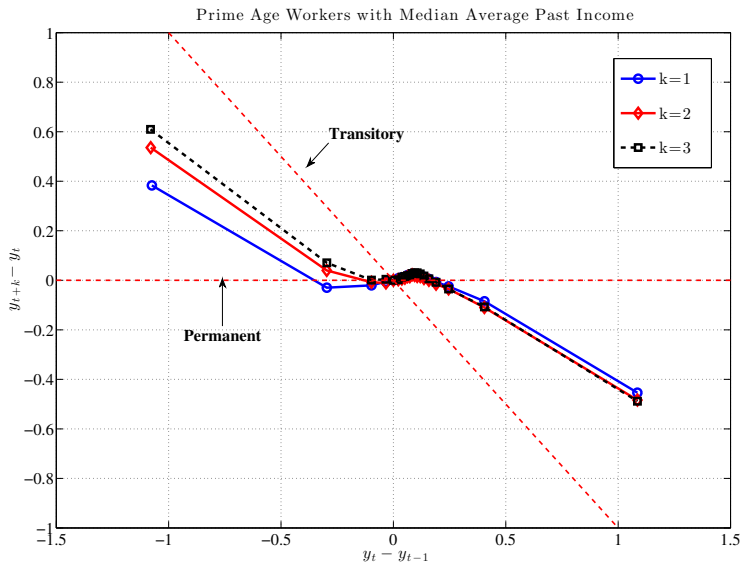
# IMPULSE RESPONSE FUNCTIONS



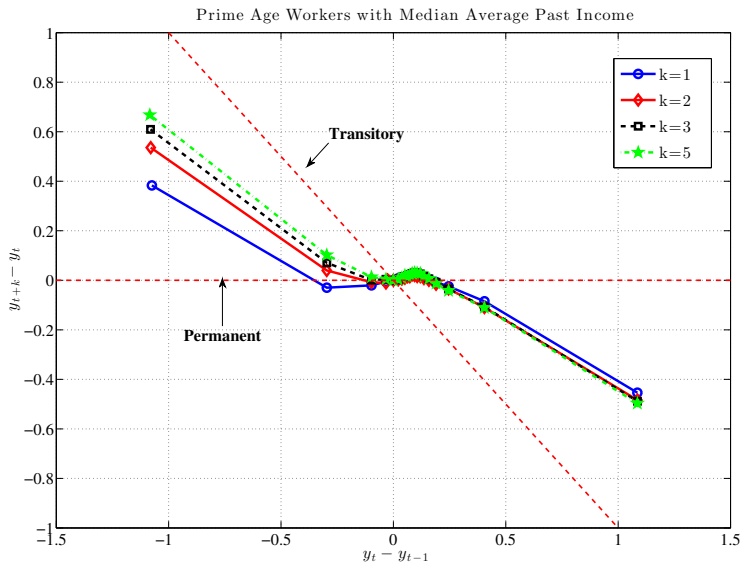
# IMPULSE RESPONSE FUNCTIONS



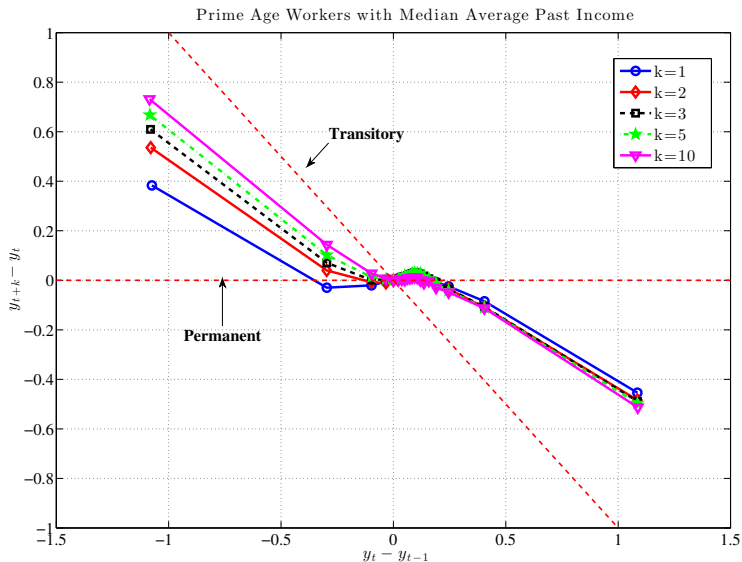
# IMPULSE RESPONSE FUNCTIONS



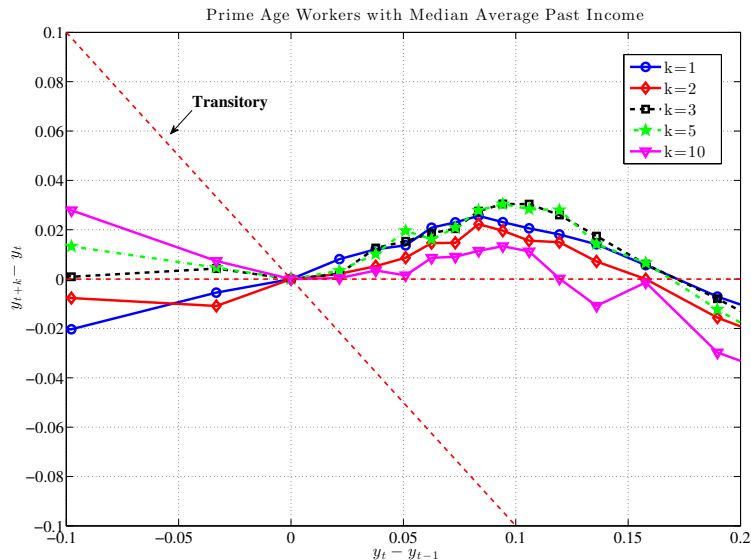
# IMPULSE RESPONSE FUNCTIONS



# IMPULSE RESPONSE FUNCTIONS

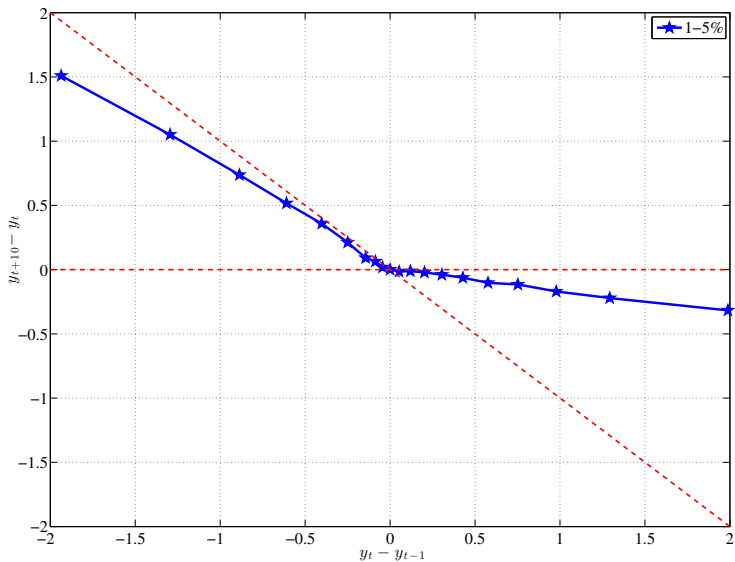


## IMPULSE RESPONSE FUNCTIONS

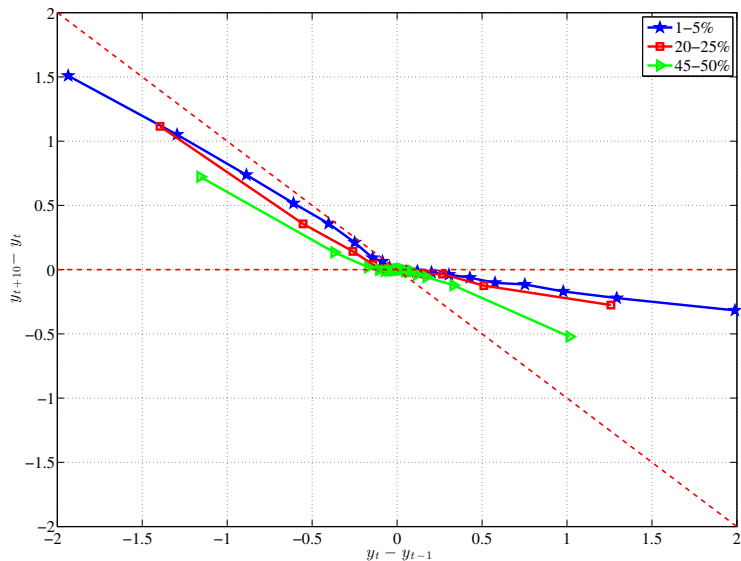




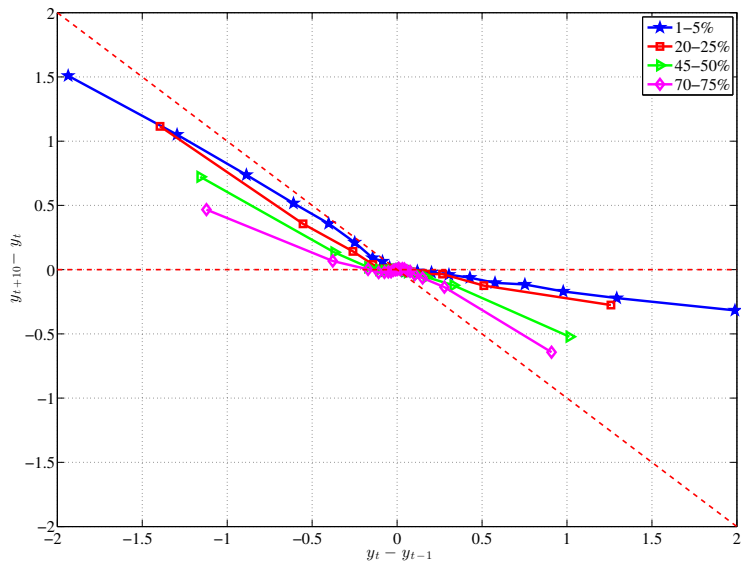
## ASYMMETRIC MEAN REVERSION



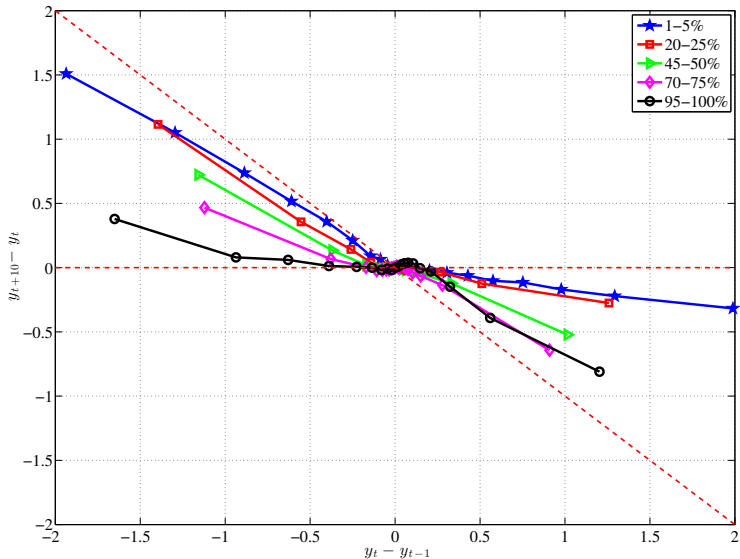
## ASYMMETRIC MEAN REVERSION



## ASYMMETRIC MEAN REVERSION



## ASYMMETRIC MEAN REVERSION



## FOUR SETS OF EMPIRICAL FACTS

1. Average income growth over the life cycle
2. Cross-sectional moments of earnings growth:  $y_{t+k} - y_t$
3. Short- and long-run dynamics of income growth
4. “Scarring” Effects of Long-Term Unemployment



ESTIMATION

## ECONOMETRIC SPECIFICATION

$$y_t^i = \underbrace{[\alpha^i + \beta^i t + \gamma^i t^2]}_{\text{HIP}} + \underbrace{z_{1,t}^i + z_{2,t}^i + z_{3,t}^i}_{\text{mixture of AR(1)s}} + \underbrace{\varepsilon_t^i}_{\text{i.i.d.}}$$

For  $j = 1, 2, 3$ :

$$z_{j,t}^i = \rho_j z_{j,t-1}^i + \eta_{j,t}^i$$



# ECONOMETRIC SPECIFICATION

$$y_t^i = \underbrace{[\alpha^i + \beta^i t + \gamma^i t^2]}_{\text{HIP}} + \underbrace{z_{1,t}^i + z_{2,t}^i + z_{3,t}^i}_{\text{mixture of AR(1)s}} + \underbrace{\varepsilon_t^i}_{\text{i.i.d.}}$$

For  $j = 1, 2, 3$ :

$$z_{j,t}^i = \rho_j z_{j,t-1}^i + \eta_{j,t}^i$$

$$\eta_{jt}^i = \begin{cases} \sim \mathcal{N}(\mu_j, \sigma_j) & \text{w.p. } \rho_j \\ 0 & \text{w.p. } 1 - \rho_j \end{cases}$$

# ECONOMETRIC SPECIFICATION

$$y_t^i = \underbrace{[\alpha^i + \beta^i t + \gamma^i t^2]}_{\text{HIP}} + \underbrace{z_{1,t}^i + z_{2,t}^i + z_{3,t}^i}_{\text{mixture of AR(1)s}} + \underbrace{\varepsilon_t^i}_{\text{i.i.d.}}$$

For  $j = 1, 2, 3$ :

$$z_{j,t}^i = \rho_j z_{j,t-1}^i + \eta_{j,t}^i$$

$$\eta_{jt}^i = \begin{cases} \sim \mathcal{N}(\mu_j, \sigma_j) & \text{w.p. } \rho_j \\ 0 & \text{w.p. } 1 - \rho_j \end{cases}$$

and

$$\begin{aligned} p_j(t, z_{t-1}) &= a_j + b_j \times z_{t-1} + c_j \times t + d_j \times z_{t-1} \times t \\ 1 &= p_1 + p_2 + p_3 \end{aligned}$$

# ESTIMATION RESULTS

<i>Specification:</i>	(2)	(4)	(8)
	<b>Benchmark model</b>	<b>Best fit</b>	<b>Standard model</b>
HIP order	1	1	0
AR(1)	2	3	1
RW	1	0	0
Heterog. variances?	yes	yes	no
<i>Parameter</i>			
$\sigma_\alpha$	0.552	0.475	0.673
$\sigma_\beta \times 10$	0.130	0.137	—
$\text{corr}_{\alpha\beta}$	-0.49	-0.04	—
<b>Objective value</b>	<b>17.17</b>	<b>15.29</b>	<b>38.75</b>

*Note: \* $\rho_\nu = 1.0$  is imposed.*

## ESTIMATION RESULTS

<i>Specification:</i>	(2)	(4)	(8)
	<b>Benchmark model</b>	<b>Best fit</b>	<b>Standard model</b>
$\rho_z$	0.259	0.083	1.00
$\rho_x$	0.425	0.512	—
$\rho_\nu$	1.0*	0.796	—
$\bar{\sigma}_z$	0.847	0.656	0.172
$\bar{\sigma}_x$	0.361	0.337	—
$\bar{\sigma}_\nu$	0.087	0.082	—
$\sigma_\varepsilon$	0.029	0.019	0.040
$\sigma_{\nu 0}$	—	0.110	—
$\sigma_{zz}$	0.107	0.167	—
$\sigma_{xx}$	0.194	0.294	—

## SELECTED ESTIMATED PARAMETERS

		$Z_1$	$Z_2$	$Z_2$
$p_j$	Probability	0.11	0.23	0.68
$\rho$	Persistence	0.08	0.51	0.80
$\sigma$	Innov. std. dev	0.66	0.34	0.08
$\sigma_\alpha$		0.48		
$\sigma_\beta \times 10$		0.14		

## WHAT TO USE IN CALIBRATION?

- These estimated processes are complex and richly parameterized.
  - How to use them for calibration?

## WHAT TO USE IN CALIBRATION?

- These estimated processes are complex and richly parameterized.
  - How to use them for calibration?
- We intend to construct Markov transition matrices that summarize these processes.

## WHAT TO USE IN CALIBRATION?

- These estimated processes are complex and richly parameterized.
  - How to use them for calibration?
- We intend to construct Markov transition matrices that summarize these processes.
- Catalan-Civale-Fazilet (2015) explore how to do this for processes with excess kurtosis and large skewness.
  - Results to so far quite encouraging.



## A STRUCTURAL INTERPRETATION

- Within-job earnings changes are small.
  - Every once in a while: find a better job or lose the job.

## A STRUCTURAL INTERPRETATION

- Within-job earnings changes are small.
  - Every once in a while: find a better job or lose the job.
- Job mobility declines with age and wage.
  - **Kurtosis** goes up with age and wage
    - **Variance** of income changes decline with age and wage

## A STRUCTURAL INTERPRETATION

- Within-job earnings changes are small.
  - Every once in a while: find a better job or lose the job.
- Job mobility declines with age and wage.
  - **Kurtosis** goes up with age and wage
    - **Variance** of income changes decline with age and wage
- **Skewness**: Job losses contribute to the left tail.
  - Larger left tail for older people and for high wage people.

## A STRUCTURAL INTERPRETATION

- Within-job earnings changes are small.
  - Every once in a while: find a better job or lose the job.
- Job mobility declines with age and wage.
  - **Kurtosis** goes up with age and wage
    - **Variance** of income changes decline with age and wage
- **Skewness**: Job losses contribute to the left tail.
  - Larger left tail for older people and for high wage people.
- These insights are mostly missed in income dynamics literature.

# CONCLUSIONS

- New empirical facts regarding individual earnings.

# CONCLUSIONS

- New empirical facts regarding individual earnings.
- Existing specifications do not capture these salient features of the data.

# CONCLUSIONS

- New empirical facts regarding individual earnings.
- Existing specifications do not capture these salient features of the data.
- We propose a richer specification that captures many of these patterns.

# CONCLUSIONS

- New empirical facts regarding individual earnings.
- Existing specifications do not capture these salient features of the data.
- We propose a richer specification that captures many of these patterns.
- Excess kurtosis and mixture of AR(1)'s can explain a well-known puzzle in the CME literature.
  - Ongoing research (Güvenen & Tonetti (2014)) shows this more rigorously.



# CONCLUSIONS

- New empirical facts regarding individual earnings.
- Existing specifications do not capture these salient features of the data.
- We propose a richer specification that captures many of these patterns.
- Excess kurtosis and mixture of AR(1)'s can explain a well-known puzzle in the CME literature.
  - Ongoing research (Guvenen & Tonetti (2014)) shows this more rigorously.
- New benchmarks and targets for calibration.

# INITIAL CONDITIONS

$$\eta_{jt}^{*i} \sim \mathcal{N}(\mu_j, \sigma_j^i) \quad \text{and} \quad \eta_{jt}^i = \eta_{jt}^{*i} \times \mathbb{I}\{\mathbf{s}_{i,t} \in I_{p_j}\} \quad (1)$$

$$\log \sigma_j^i \sim \mathcal{N}(\bar{\sigma}_j - \frac{\sigma_{jj}^2}{2}, \sigma_{jj}^2), j = z, x, \quad \sigma_v^i \equiv \sigma_v \quad (2)$$