# Sources of Increasing Earnings Inequality: Reconciling Survey and Administrative Data

John Haltiwanger[1]    Henry Hyatt[2]    James Spletzer[3]

[1]Department of Economics, University of Maryland

[2]Center for Economic Studies, U.S. Census Bureau

[3]Retired

NBER Conference on Wage Dynamics in the 21st Century

September 16, 2022

We compare two very different perspectives on rising earnings inequality.

A large literature using household survey data emphasizes:

▶ rising dispersion across education and occupation groups.

▶ industry effects are modest or offsetting.

A more recent literature using matched employer-employee admin data emphasizes:

▶ rising dispersion between firms (and industries)

▶ rising between firm and industry dispersion is accounted for by sorting and segregation.

Our analysis uses linked micro admin and survey data to reconcile these approaches.

# Industries and increasing inequality

What is the role of industry in increasing inequality?

**Administrative records data**

*Most of the rise in overall earnings inequality is accounted for by rising between-industry inequality.*

- ▶ Haltiwanger, Hyatt, and Spletzer (2022)

*The dominant driver of the rising inequality of both earnings and wage rates in Italy is the growing heterogeneity of pay across industries.*

- ▶ Briskar, Di Porto, Rodriguez Mora, and Tealdi (2022)

**Survey data**

*The between-group variance component linked to industry has been declining over time.*

- ▶ Hoffmann, Lee, and Lemieux (2020)

*Using the CPS, we show that since the 1980s there has been a decline of about one third in the dispersion of industry wage premia.*

- ▶ Stansbury and Summers (2020)

- From -5.9% to 61.9%

- Hoffmann, Lee, and Lemieux (2020, HLL) use data from the Annual Social and Economic Supplement of the Current Population Survey (CPS)

- Haltiwanger, Hyatt, and Spletzer (2022, HHS) use Longitudinal Employer- Household Dynamics (LEHD) administrative records data for 18 U.S. states

- We link these datasets (later)

- There are methodological differences (later)

- ▶ We conducted a number of exercises to compare HLL with our findings in HHS.

- ▶ Our first exercise concerns the industry classification method.

- ▶ In the late 1990s, the North American Industrial Classification System (NAICS) replaced the Standard Industrial Classification (SIC)

- ▶ Using HLL's method (next slide) and 18 NAICS sectors yields an industry contribution of 0.8% rather than the -5.9% when using 12 SIC aggregate industries

**Our replication of HLL:** For the seven-year intervals $\{1996\text{-}02, 2012\text{-}18\}$, we estimate a human capital earnings ($y_i$ for worker $i$) equation in three steps:

$$y_i = AgeEduc_i\beta_1 + \varepsilon_i \tag{1}$$

$$y_i = AgeEduc_i\beta_1 + Occupation_i\beta_2 + \varepsilon_i \tag{2}$$

$$y_i = AgeEduc_i\beta_1 + Occupation_i\beta_2 + Industry_i\beta_3 + \varepsilon_i \tag{3}$$

Estimating equation (1) provides the $R^2$ from including age and education alone. The contribution of occupation is the $R^2$ from equation (2) minus the $R^2$ from (1). The contribution of industry is the $R^2$ from equation (3) minus the $R^2$ from (2).

**Our between-industry in HHS:** We use the simple decomposition

$$\underbrace{var(y_{i,k} - \bar{y})}_{\substack{\text{earnings} \\ \text{variance}}} = \underbrace{var(y_{i,k} - \bar{y}_k)}_{\substack{\text{within-industry} \\ \text{dispersion}}} + \underbrace{var(\bar{y}_k - \bar{y})}_{\substack{\text{between-industry} \\ \text{dispersion}}} \tag{4}$$

where $\bar{y}_k$ is average earnings in industry $k$

- ▶ Using HLL's estimation method and 18 NAICS industries yields an industry contribution (pay premium) of 0.8%

- ▶ In this same CPS dataset, the between-industry component explains 23.1% of the increase

- ▶ The difference between these estimates reflects the way workers sort and segregate into different industries on the basis of age, education, and occupation (more about this later)

**HLL**: annual real earnings > \$7840, and:

- ► weeks worked > 49
- ► usual hours ≥ 40
- ► real hourly wage > \$4

**HHS**: annual earnings > \$3770

There are additional differences in which ages, job types, etc. are included.

**Common coding applied to both datasets**

Annual real earnings > $3770

Other common coding:

- ▶ Exclude self-employed
- ▶ Age 20-60
- ▶ Exclude longest job if government
- ▶ Use PCE (2013=100)
- ▶ Any firm size

- The between-industry contribution is 29.3% after common coding

- As opposed to 23.1% when we follow HLL's sample selection method

The bar chart shows "Industry share of increasing inequality 1996-02 to 2012-18" (y-axis) with categories: Our HLL replication, 18 NAICS sectors, CPS between-industry, Common coding, Linked CPS-LEHD, Our estimate in HHS.

- ▶ After linking our common coded CPS to the LEHD, the between-industry contribution rises from 29.3% to 46.0%

- ▶ This jump is driven by 18 vs 50 state differences that are much larger than those found in the LBD or QCEW.

- ▶ The 18 vs. 50 state differences the CPS reports in the Retail Trade and Information sectors are much larger than those obtained from published aggregates.

- ▶ Within- vs. between patterns are similar in 50 vs. 18 states in administrative data.

In our linked CPS-LEHD dataset, at the NAICS sector level, there is less than 50% agreement in:

- ▶ Wholesale Trade (46.9%)
- ▶ Information (46.6%)
- ▶ Educational Services (44.6%)
- ▶ Other Services (48%)

- ▶ Now that we have linked the CPS to the LEHD, we can use employer-reported industries
- ▶ The between-industry component explains 52.2% of the increase using 18 NAICS sectors from LEHD
    - ▶ as opposed to 46.0% when using 18 NAICS sectors from household-reported CPS industries
- ▶ Using 299 4-digit NAICS industry groups brings us to 65.5%
- ▶ We reported 61.9% in HHS

This is similar when we replace CPS earnings with LEHD earnings (66.2%), and when we use the full LEHD with common coding applied (64.5%)

Taking stock (-5.9% vs. 64.5%):

▶ Industry pay premium vs. between-industry variation (31% of increase)

▶ Differences in the unlinked versus linked CPS-LEHD samples (23% of increase)

▶ Using LEHD 4-digit NAICS industry rather than CPS NAICS sectors (27% of increase)

Pay premia vs. between-industry differences: the roles of sorting and segregation

Song, Price, Guvenen, Bloom, and von Wachter (2019) build on Abowd, Kramarz, and Margolis (1999) and Card, Heining, and Kline (2013) to measure between-firm inequality in terms of the following:

▶ **Pay premia**: some firms offer greater earnings to any worker

▶ **Sorting:** high-paying firms employ more highly paid workers

▶ **Segregation:** highly paid workers concentrate among each other

In HHS, we extend the Song et al. (2019) framework in order to measure how these components of inequality occur between vs. within industries

In this paper, we show how to apply this method to the CPS to measure industry-level sorting and segregation by age, education, and occupation

**Pay premia vs. between-industry:** We re-write the human capital earnings equation used by HLL (introducing a subscript for industry $k$) as

$$y_{i,k} = Z_{i,k}\beta_Z + Industry_{i,k}\beta_3 + \varepsilon_{i,k}, \tag{5}$$

where $Z$ concatenates the $AgeEduc_i$ and $Occupation_i$ vectors, and $\beta_Z$ concatenates the marginal effects vectors $\beta_1$ and $\beta_2$.

Define $\overline{Z_k\beta_Z}$ as the industry mean of $Z_{i,k}\beta_Z$.

Taking variances of both sides of the human capital earnings equation results in:

$$\underbrace{var(y_{i,k})}_{\substack{\text{earnings} \\ \text{variance}}} = \underbrace{var(Z_{i,k}\beta_Z - \overline{Z_k\beta_Z})}_{\substack{\text{within-industry dispersion} \\ \text{from age, education,} \\ \text{and occupation}}} + \underbrace{var(\overline{Z_k\beta_Z})}_{\substack{\text{between-industry} \\ \text{segregation}}} +$$

$$\underbrace{var(Industry_{i,k}\beta_3)}_{\substack{\text{between-industry} \\ \text{pay premium}}} + \underbrace{2cov(\overline{Z_k\beta_Z}, Industry_{i,k}\beta_3)}_{\text{between-industry sorting}} + \underbrace{var(\varepsilon_{i,k})}_{\substack{\text{residual dispersion} \\ \text{(within-industry)}}} \tag{6}$$

## Using the CPS alone

- ▶ We start by applying our decomposition to our replication of the HLL analysis dataset

- ▶ Pay premium explains -6.9% of the increase using 12 SIC aggregates, 1.0% using 18 NAICS sectors

- ▶ Segregation explains 13.3%-14.8%

- ▶ Sorting explains 7.3%-7.4%

- ▶ Within-industry dispersion by age, education, & occupation explains 13.3% (18.2%) using 12 SIC aggregates (18 NAICS sectors).

**Using the linked CPS-LEHD**

- ▶ 299 4-digit NAICS industries, common coded

- ▶ Still using age, education, and occupation following equation (6)

- ▶ Industry dispersion contributes -1.2% using CPS earnings

- ▶ Industry dispersion contributes 22.0% using LEHD earnings

- ▶ Segregation 35.1%, sorting 31.7% using CPS earnings

- ▶ Segregation 29.4%, sorting 14.8% using LEHD earnings

**If we estimate AKM and implement our decomposition from HHS**

- ▶ Using person and firm effects rather than age, education, & occupation
- ▶ Pay premium explains 9.3%-9.6%
- ▶ Segregation explains 29.2%-29.6%
- ▶ Sorting explains 25.9%-26.9%
- ▶ Within-industry worker and firm effects explain 37.2%-40.5%
- ▶ Small offsetting effect of residual

# Industries and occupations

We use public domain data from the Occupational Employment and Wage Statistics (OEWS) to construct a dataset of 287 4-digit NAICS industries $k$ by 22 occupations $j$.

We estimate the following equation for the intervals 2002-03 and 2015-16:

$$y_{j,k} = Occupation_{j,k}\beta_2 + Industry_{j,k}\beta_3 + \varepsilon_{j,k}. \tag{7}$$

Taking (employment-weighted) variances of both sides of equation (7) yields:

$$\underbrace{var(y_{j,k})}_{\substack{\text{earnings} \\ \text{variance}}} = \underbrace{var(Occupation_{j,k}\beta_2 - \overline{Occupation_k\beta_2})}_{\substack{\text{within-industry dispersion} \\ \text{from occupation}}} +$$

$$\underbrace{var(\overline{Occupation_k\beta_2})}_{\substack{\text{between-industry} \\ \text{segregation}}} + \underbrace{var(Industry_{j,k}\beta_3)}_{\substack{\text{between-industry} \\ \text{pay premia}}} +$$

$$\underbrace{2cov(\overline{Occupation_k\beta_2}, Industry_{j,k}\beta_3)}_{\text{between-industry sorting}} + \underbrace{var(\varepsilon_{j,k})}_{\substack{\text{residual dispersion} \\ \text{(within-industry)}}} \tag{8}$$

## Using OEWS aggregates



- ▶ 287 NAICS industries

- ▶ Pay premium explains 2.8%

- ▶ Segregation explains 31.3%

- ▶ Sorting explains 53.7%

- ▶ Within-industry occupation effects explain 16.0%

- ▶ Small offsetting effect of residual

- ▶ Between-industry occupation changes also dominate within-industry in the CPS

- ▶ Sorting and segregation reflect changes in the way workers are allocated across industries

# Conclusion

Findings from the CPS, OEWS, and our administrative records data:

1. A large share of inequality growth in recent decades has occurred at the industry level

2. The role of industry pay premia in increasing inequality is relatively small

3. Sorting and segregation are of first order importance when assessing the role of industries in inequality growth
   ▶ whether we measure sorting and segregation using education, occupation, or AKM worker effects

# Appendix

Our 18 LEHD states are: CA, CO, CT, HI, ID, IL, KS, LA, MD, MN, MT, NC, NJ, OR, RI, TX, WA, and WI. Back: introduction Back: linked datasets

| Criterion | HLL CPS-ASEC | Common Coding | HHS LEHD |
|-----------|--------------|---------------|----------|
| Earnings | Wage & Salary + Self Employment + Farm | Wage & Salary | Wage & Salary |
| Age | 26-65 | 20-60 | 20-60 |
| Top coding | Truncate top 1% each year (by gender) | Mean of top 0.001% pooled all years | Mean of top 0.001% pooled all years |
| Bottom coding | Weeks worked > 49 & usual hours > 40 & real hourly wage > $4 & annual real earnings > $7840 | Annual real earnings > $3770 | Annual real earnings > $3770 |
| Government jobs | Include all government jobs | Exclude longest job last year that is government | Exclude all government jobs |
| Deflator | CPI (2018=100) | PCE (2013=100) | PCE (2013=100) |
| Firm size | Any | Any | Firm Size > 20 |

Recall our formula for between- vs. within-industry variance:

$$\underbrace{\text{var}(y_{i,k} - \bar{y})}_{\substack{\text{earnings} \\ \text{variance}}} = \underbrace{\text{var}(y_{i,k} - \bar{y}_k)}_{\substack{\text{within-industry} \\ \text{dispersion}}} + \underbrace{\text{var}(\bar{y}_k - \bar{y})}_{\substack{\text{between-industry} \\ \text{dispersion}}}$$

The between-industry component can be computed for published aggregates of industry-level average earnings. We use the Quarterly Census of Employment and Wages (QCEW). The 50 state between-industry variance growth is lower than that of our 18 states in both the CPS and the QCEW, but this ratio is lower in the CPS (70.5%) than in the QCEW (79.6%).

$$\underbrace{\text{var}(y_{i,k} - \bar{y})}_{\substack{\text{earnings} \\ \text{variance}}} = \underbrace{\text{var}(y_{i,k} - \bar{y}_k)}_{\substack{\text{within-industry} \\ \text{dispersion}}} + \underbrace{\text{var}(\bar{y}_k - \bar{y})}_{\substack{\text{between-industry} \\ \text{dispersion}}}$$

Differences between the CPS and QCEW are driven by two industries:

| Data | CPS (Micro) | | CPS (Agg) | | QCEW (Agg) | |
|------|-----------|-----------|-----------|-----------|-----------|-----------|
| | 50 states | 18 states | 50 states | 18 States | 50 States | 18 States |
| Contribution to variance growth from 1996-02 to 2012-18: | | | | | | |
| Retail Trade | 0.0035 | 0.0054 | 0.0024 | 0.0041 | 0.0066 | 0.0078 |
| Information | 0.0054 | 0.0082 | 0.0050 | 0.0076 | 0.0030 | 0.0038 |
| | | | | | | |
| Ratio of 50 State to 18 State (%): | | | | | | |
| Retail Trade | 64.8% | | 59.3% | | 84.8% | |
| Information | 65.9% | | 65.8% | | 78.7% | |

Use the Longitudinal Business Database to measure the share of inequality growth that is within-industry

Average earnings at the establishment- or EIN-level, so omits variation between workers at the same employer

18 states have an industry share that is 0.8 to 3.4 percentage points higher than the national average

Back

- ▶ Replacing CPS earnings with LEHD earnings yields a similar between-industry share of the increase in earnings dispersion (66.2% vs. 65.5%)

- ▶ This is despite the higher variance of LEHD earnings (relative to CPS earnings), even after common coding and linking

- Applying common coding to the LEHD earnings in the full LEHD dataset brings us to 64.5%

- We reported 61.9% in HHS

Following Abowd, Kramarz, and Margolis (1999, AKM), assume that earnings $y_t^{i,j,k,p}$ is the sum of the worker $i$ effect $\theta^{i,p}$, a firm $j$ in industry $k$ effect $\psi^{j,k,p}$, and observable characteristics $X_t^{i,p}$ (marginal effects $\beta^p$).

$$y_t^{i,j,k,p} = X_t^{i,p}\beta^p + \theta^{i,p} + \psi^{j,k,p} + \varepsilon_t^{i,j,k,p} \qquad (9)$$

We estimate this AKM equation separately by interval $p$.

For the purposes of this presentation (on following slide), we omit the superscript $p$ and the effects of observable characteristics $X_t^{i,p}\beta^p$.

Following Song et al. (2019), the variance of earnings can be written in terms of firm-level averages $\bar{\theta}^{j,k}$ (worker effects)

$$\text{var}(y_t^{i,j,k}) = \underbrace{\text{var}(\bar{\theta}^{j,k})}_{\substack{\text{total} \\ \text{segregation}}} + \underbrace{\text{var}(\psi^{j,k})}_{\substack{\text{total pay} \\ \text{premia}}} + \underbrace{2\text{cov}(\bar{\theta}^{j,k}, \psi^{j,k})}_{\text{total sorting}} +$$
$$\underbrace{\text{var}(\theta^i - \bar{\theta}^{j,k})}_{\text{within-firm person effects}} + \underbrace{\text{var}(\varepsilon_t^{i,j,k})}_{\text{residual}} \tag{10}$$

## Comparison to Song et al. (2019), males only

| | Song et al. (2019) growth 1994-2000 2007-2013 | LEHD growth 1996-2002 2012-2018 |
|---|---|---|
| Total variance increase | 0.096 | 0.126 |
| Within-firm share | 13.5% | 15.5% |
| Between-firm share | 86.5% | 84.5% |
| Segregation | 35.5% | 37.4% |
| Pay premia | 14.6% | 11.8% |
| Sorting | 37.5% | 35.3% |

*Notes*: Estimates from Table V (page 36) of Song et al. (2019). LEHD estimates from Haltiwanger, Hyatt, and Spletzer (2022).

Following Song et al. (2019), the variance of earnings can be written in terms of firm-level averages $\bar{\theta}^{j,k}$ (worker effects). In HHS, we introduce industry-average worker effects $\bar{\theta}^k$ and firm effects $\bar{\psi}^k$. The variance of earnings can now be written:

$$
\begin{aligned}
\mathrm{var}(y_t^{i,j,k}) = \underbrace{\mathrm{var}(\bar{\theta}^k)}_{\substack{\text{between-industry}\\\text{segregation}}} \;+\; \underbrace{\mathrm{var}(\bar{\theta}^{j,k} - \bar{\theta}^k)}_{\substack{\text{within-industry, between-firm}\\\text{segregation}}} \;+\; \\
\underbrace{\mathrm{var}(\bar{\psi}^k)}_{\substack{\text{between-industry}\\\text{pay premia}}} \;+\; \underbrace{\mathrm{var}(\psi^{j,k} - \bar{\psi}^k)}_{\substack{\text{within-industry, between-firm}\\\text{pay premia}}} \;+\; \\
\underbrace{2\mathrm{cov}(\bar{\theta}^k, \bar{\psi}^k)}_{\text{between-industry sorting}} \;+\; \underbrace{2\mathrm{cov}[(\bar{\theta}^{j,k} - \bar{\theta}^k), (\psi^{j,k} - \bar{\psi}^k)]}_{\text{within-industry, between-firm sorting}} + \\
\underbrace{\mathrm{var}(\theta^i - \bar{\theta}^{j,k})}_{\text{within-firm person effects}} \;+\; \underbrace{\mathrm{var}(\varepsilon_t^{i,j,k})}_{\text{residual}}
\end{aligned}
\tag{11}
$$

Back: sorting and segregation Back: empirical results

| Data | CPS | Linked CPS-LEHD | Linked CPS-LEHD |
|---|---|---|---|
| Sample | HLL JEP | Common coded | Common coded |
| Earnings measure | CPS | CPS | LEHD |
| Industry measure | CPS 18 | LEHD 299 | LEHD 299 |
| **Within-industry:** | | | |
| | | | |
| Age, education & occupation: | 18.2% | 6.0% | 13.3% |
| Age and education | 11.8% | 9.0% | 14.4% |
| **Occupation** | **4.1%** | **0.4%** | **0.0%** |
| **Covariance: age+educ & occ.** | **2.2%** | **-3.6%** | **-1.0%** |
| | | | |
| Residual | 58.8% | 28.5% | 20.6% |
| | | | |
| **Between-industry:** | 23.1% | 65.5% | 66.2% |
| | | | |
| Segregation | 14.8% | 31.7% | 14.8% |
| Age and education | 3.8% | 9.6% | 5.6% |
| **Occupation** | **2.5%** | **5.4%** | **2.0%** |
| **Covariance: age+educ. & occ:** | **8.4%** | **16.7%** | **7.0%** |
| | | | |
| Pay premia | 1.0% | -1.2% | 22.0% |
| | | | |
| Sorting | 7.3% | 35.1% | 29.4% |
| Covariance: age+educ. & ind. | 2.0% | 19.9% | 19.0% |
| **Covariance: industry & occ.** | **5.3%** | **15.5%** | **10.5%** |

Change from 1996-02 to 2012-18

Within-industry occupation dispersion explains **-3.2%** to **6.3%** of the increase in earnings dispersion.

Between-industry sorting and segregation by occupation explains **16.2%** to **37.6%** of increasing dispersion

Most of the contribution of occupations to increasing inequality is through how they are allocated across industries

# The top thirty industries

- **High-tech**: 11 of the 19 high-paying industries are high-tech in terms of STEM intensity as classified by Hecker (2005) and Goldschlag and Miranda (2016)
  - One-third of the increase in between-industry inequality

- **Mining**: 2 high-paying: oil and gas (also high-tech), drilling wells

- **Finance and Insurance**: 4 of the 19 high-paying industries

- **Management of Companies and Enterprises** : corporate headquarters

- **Health Care and Social Assistance**: 2 high-paying (physician offices, hospitals), 3 low-paying (in-home care, nursing homes, social services)

- **Support services**: 2 of the 11 low-paying industries

- **Retail, restaurants, and gyms**: 6 of the 11 low-paying industries
  - Another one-third of the increase in between-industry inequality

# High-tech: 11 of the 19 high-paying industries (I)

| Industry title | Employment share: average | Employment share: change | Relative log earnings: average | Relative log earnings: change | Share of bet.-ind. var. gth. |
|---|---|---|---|---|---|
| Oil & Gas Extraction | 0.3% | -0.1% | 1.012 | 0.247 | 1.8% |
| Pharmaceutical Manufacturing | 0.5% | -0.1% | 0.799 | 0.203 | 1.6% |
| Semiconductor Manufacturing | 0.8% | -0.5% | 0.556 | 0.299 | 1.4% |
| Professional Equip. Wholesaler | 0.7% | -0.0% | 0.557 | 0.190 | 1.9% |
| Software Publishers | 0.5% | 0.2% | 1.009 | 0.186 | 5.6% |
| Data Processing Services | 0.3% | -0.0% | 0.545 | 0.301 | 1.3% |

*Notes*: Persons with annual real earnings > $3770 in EINs with 20 or more employees. Average log earnings for industry $k$ are relative to the economy average. The 1996-2002 and 2012-2018 intervals are averaged. Changes are the growth (or decline) from 1996-2002 to 2012-2018.

Back

| Industry title | Employment share: | | Relative log earnings: | | Share of bet.-ind. var. gth. |
|---|---|---|---|---|---|
| | average | change | average | change | |
| Other Information Services | 0.2% | 0.3% | 0.798 | 0.699 | 5.8% |
| Architectur. & Enginr. Services | 1.2% | 0.1% | 0.469 | 0.161 | 2.6% |
| Computer Systems Design | 1.7% | 0.9% | 0.663 | 0.012 | 5.6% |
| Management & Scientific Serv. | 0.9% | 0.6% | 0.381 | 0.069 | 1.8% |
| Scientific Research Services | 0.8% | -0.1% | 0.741 | 0.244 | 3.3% |

*Notes*: Persons with annual real earnings > $3770 in EINs with 20 or more employees. Average log earnings for industry $k$ are relative to the economy average. The 1996-2002 and 2012-2018 intervals are averaged. Changes are the growth (or decline) from 1996-2002 to 2012-2018.

Back

# Mining: 2 of the 19 high-paying industries

| Industry title | Employment share: average | Employment share: change | Relative log earnings: average | Relative log earnings: change | Share of bet.-ind. var. gth. |
|---|---|---|---|---|---|
| Oil & Gas Extraction | 0.3% | -0.1% | 1.012 | 0.247 | 1.8% |
| Support Activities for Mining | 0.5% | 0.3% | 0.374 | 0.191 | 1.4% |

*Notes*: Persons with annual real earnings > $3770 in EINs with 20 or more employees. Average log earnings for industry $k$ are relative to the economy average. The 1996-2002 and 2012-2018 intervals are averaged. Changes are the growth (or decline) from 1996-2002 to 2012-2018.

# Finance and Insurance (NAICS sector 52), Management of Companies and Enterprises (NAICS sector 55): 5 of 19 high-paying

| Industry title | Employment share: average | change | Relative log earnings: average | change | Share of bet.-ind. var. gth. |
|---|---|---|---|---|---|
| Depository Credit Intermediat. | 2.1% | 0.0% | 0.189 | 0.234 | 2.5% |
| Securities Brokerage | 0.5% | -0.1% | 0.866 | 0.204 | 1.1% |
| Other Financial Invest. Activity | 0.3% | 0.1% | 0.834 | 0.388 | 3.3% |
| Insurance Carriers | 1.6% | -0.4% | 0.488 | 0.167 | 2.3% |
| Management of Companies | 2.0% | -0.1% | 0.471 | 0.201 | 5.0% |

*Notes*: Persons with annual real earnings > $3770 in EINs with 20 or more employees. Average log earnings for industry $k$ are relative to the economy average. The 1996-2002 and 2012-2018 intervals are averaged. Changes are the growth (or decline) from 1996-2002 to 2012-2018.

# Health Care and Social Assistance (NAICS sector 62): 2 high-paying, 3 low-paying

| Industry title | Employment share: average | change | Relative log earnings: average | change | Share of bet.-ind. var. gth. |
|---|---|---|---|---|---|
| Offices of Physicians | 1.7% | 0.5% | 0.254 | 0.099 | 1.6% |
| Home Health Care Services | 0.8% | 0.4% | -0.525 | -0.016 | 1.7% |
| General Medical & Hospitals | 4.5% | 0.5% | 0.205 | 0.162 | 4.2% |
| Continuing Care Retirement | 0.6% | 0.4% | -0.493 | -0.001 | 1.2% |
| Individual & Family Services | 0.8% | 0.6% | -0.490 | -0.155 | 3.5% |

*Notes*: Persons with annual real earnings > $3770 in EINs with 20 or more employees. Average log earnings for industry $k$ are relative to the economy average. The 1996-2002 and 2012-2018 intervals are averaged. Changes are the growth (or decline) from 1996-2002 to 2012-2018.

# Support services

| Industry title | Employment share: average | change | Relative log earnings: average | change | Share of bet.-ind. var. gth. |
|---|---|---|---|---|---|
| Employment Services | 3.9% | 0.6% | -0.685 | 0.017 | 2.5% |
| Services to Buildings & Dwell. | 1.1% | 0.3% | -0.493 | -0.002 | 1.1% |

*Notes*: Persons with annual real earnings > $3770 in EINs with 20 or more employees. Average log earnings for industry *k* are relative to the economy average. The 1996-2002 and 2012-2018 intervals are averaged. Changes are the growth (or decline) from 1996-2002 to 2012-2018.

## Retail, restaurants, and gyms

| Industry title | Employment share: average | change | Relative log earnings: average | change | Share of bet.-ind. var. gth. |
|---|---|---|---|---|---|
| Building Material & Supplies | 0.9% | 0.1% | -0.293 | -0.180 | 1.5% |
| Grocery Stores | 2.3% | 0.0% | -0.378 | -0.194 | 4.7% |
| Clothing Stores | 0.7% | -0.0% | -0.607 | -0.244 | 2.6% |
| Othr. Genrl. Merchandise Stores | 1.4% | 1.5% | -0.539 | -0.051 | 6.8% |
| Othr. Amusement & Recreation | 0.6% | 0.1% | -0.594 | -0.106 | 1.7% |
| Restaurants & Othr. Eat Places | 4.9% | 2.0% | -0.739 | -0.027 | 16.9% |

*Notes*: Persons with annual real earnings > \$3770 in EINs with 20 or more employees. Average log earnings for industry $k$ are relative to the economy average. The 1996-2002 and 2012-2018 intervals are averaged. Changes are the growth (or decline) from 1996-2002 to 2012-2018.

Back