

# 个人收入决定因素研究：基于CFPS数据的实证分析

好的，这是本研究的摘要与引言：

## 摘要

本研究旨在探究中国个体税后工资性收入的主要决定因素。基于中国家庭追踪调查（CFPS）2022年数据，采用扩展的明瑟方程，并通过加权最小二乘法（WLS）控制省份固定效应进行分析。研究发现，受教育年限、工作经验（年龄）、性别及城乡属性均对个人收入有显著影响。教育回报率约为8.5%，经验呈现倒U型，且存在显著的性别与城乡收入差距。本研究为理解当前中国收入分配格局提供了实证依据。

## 引言

个体收入差异及其决定因素是社会经济研究的核心议题。在中国经济转型背景下，理解工资性收入的形成机制对制定有效的收入分配政策至关重要。本研究利用CFPS2022截面数据，构建计量模型，重点考察教育、经验、性别及城乡等因素对个人税后工资性收入的影响。研究旨在揭示人力资本与结构性因素在塑造收入不平等中的作用，为促进社会公平提供参考。

## 研究计划

核心研究问题：本研究旨在探究中国个体收入的主要决定因素，特别是教育、工作经验（以年龄代理）、性别及城乡属性对个人税后工资性收入的量化影响。

计量经济模型：主要采用扩展的明瑟收入方程，通过普通最小二乘法（OLS）进行估计。因变量（工资性收入）将进行对数转换，以缓解数据偏度并使系数易于解释为半弹性。模型将包含年龄的线性项和平方项，以捕捉工作经验收益的非线性特征。

关键变量定义与角色：

i. 因变量: `emp_income`（过去12个月税后工资性收入），取对数后使用。选择此变量因为它直接衡量了个人的劳动市场回报。

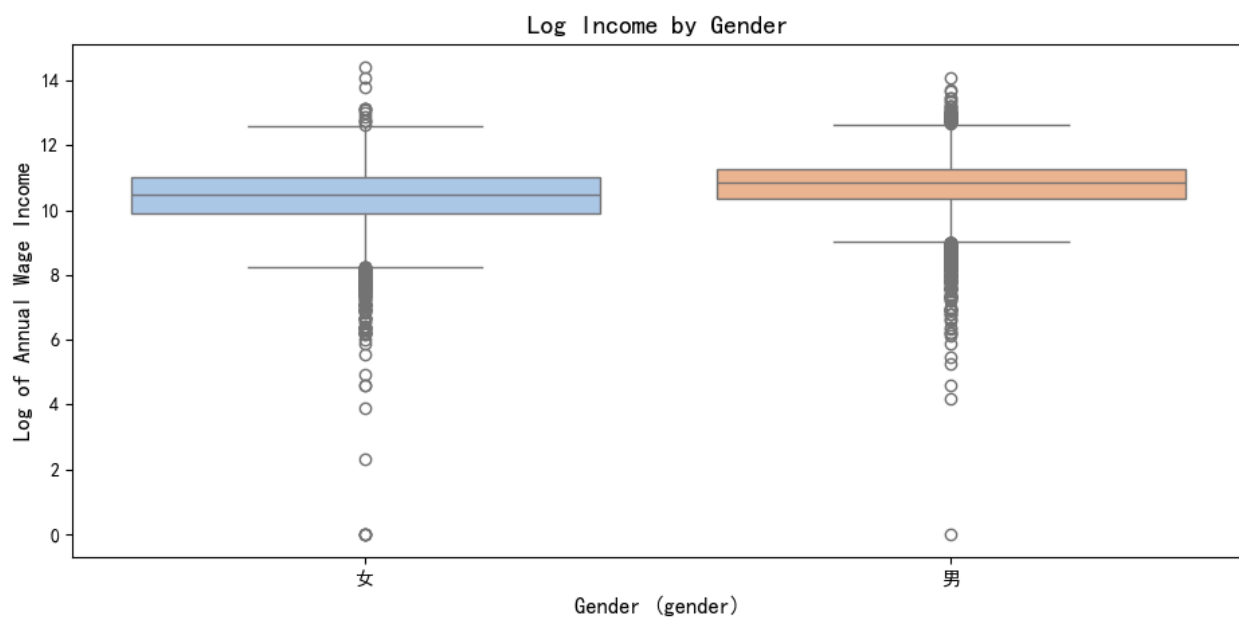
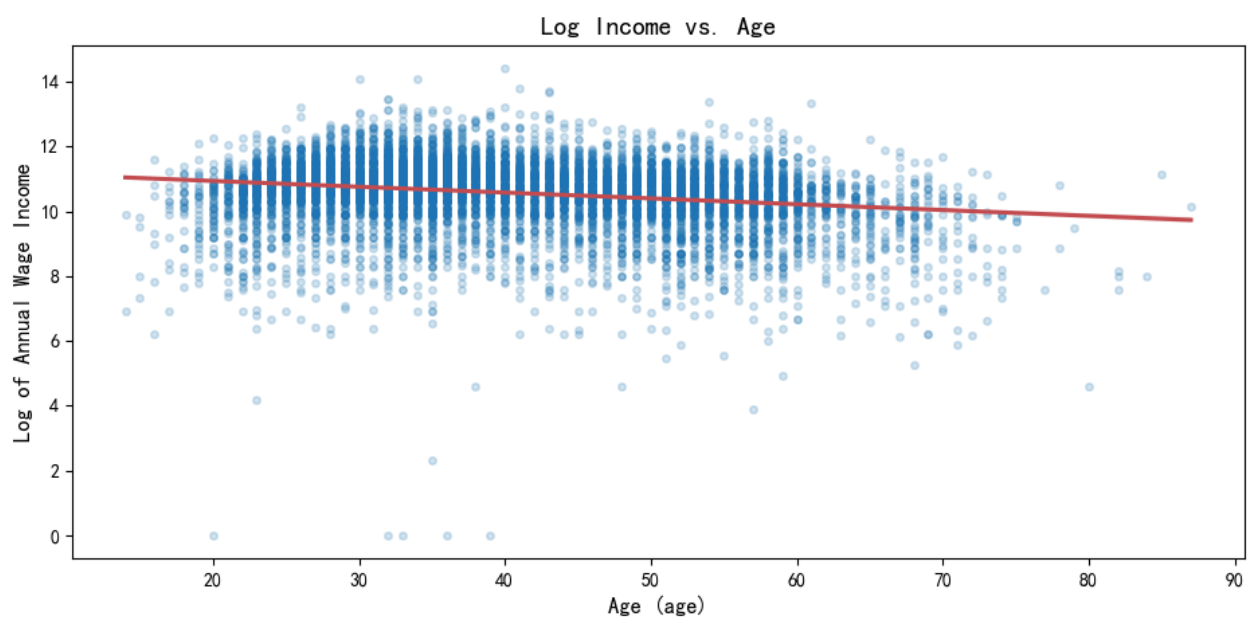
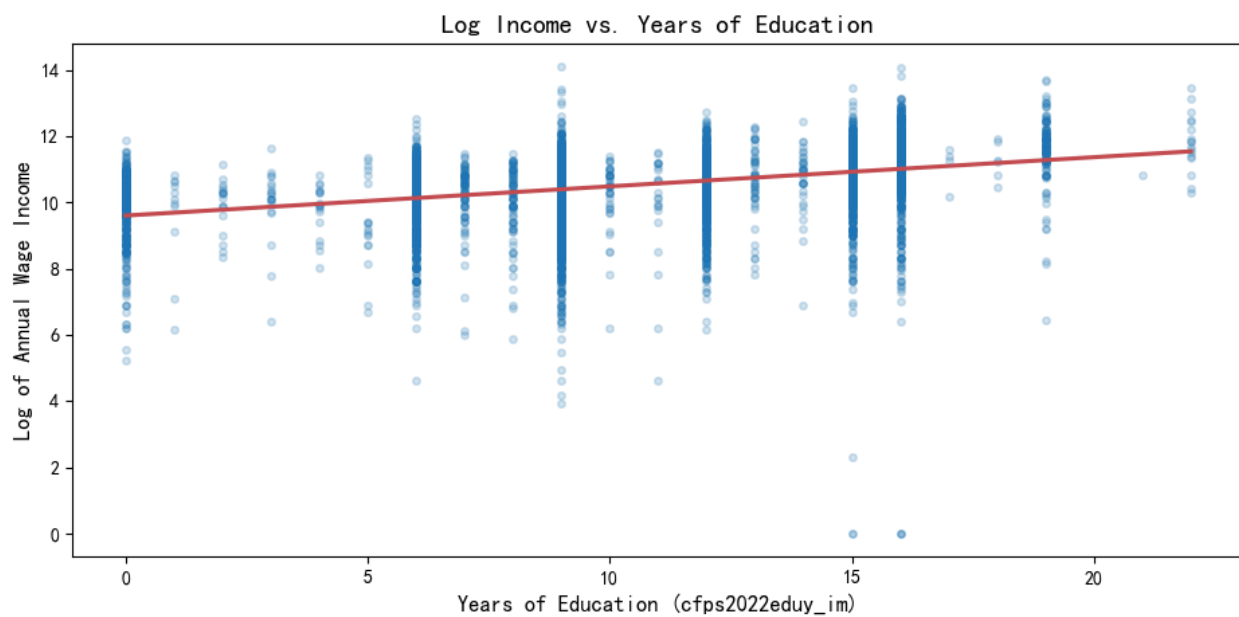
ii. 核心自变量:

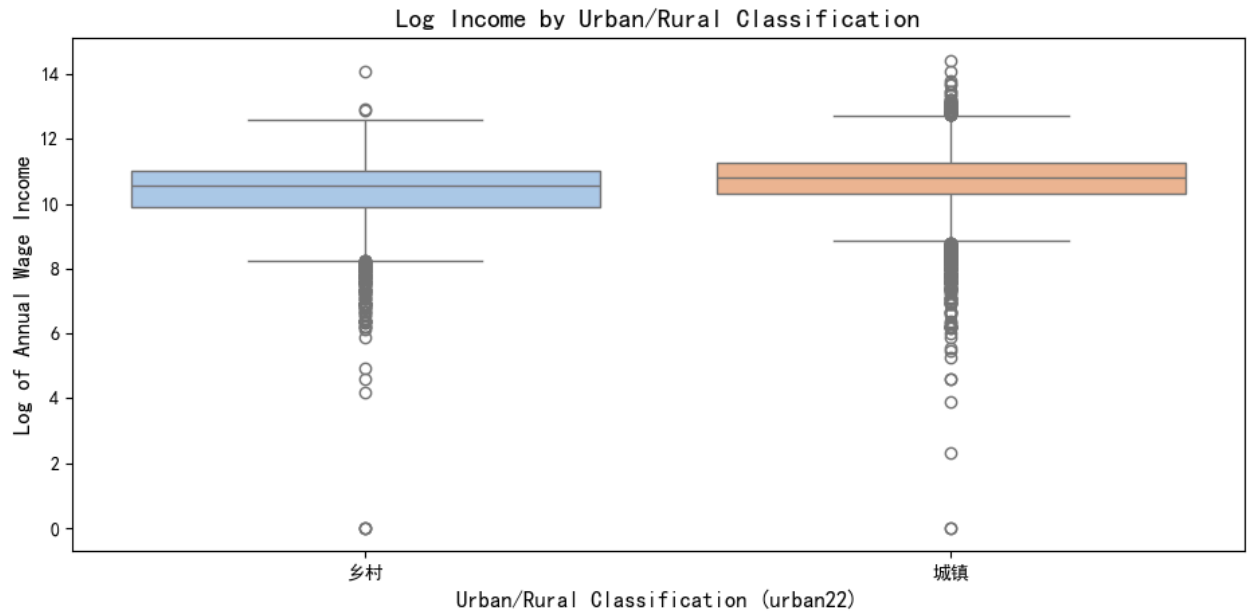
- `cfps2022eduy_im` (受教育年限): 人力资本理论核心变量，预期对收入有显著正向影响。
- `age` (年龄): 及其平方项，作为工作经验的代理，预期呈现先上升后下降的倒U型关系。
- `gender` (性别): 用于检验性别工资差异。
- `urban22` (城乡分类): 预期城镇居民收入显著高于农村居民。

iii. 控制变量:

- `qea0` (婚姻状况): 不同婚姻状态可能影响收入。
- `party` (是否党员): 政治资本可能影响收入。
- `qp201` (健康状况): 健康是重要的人力资本组成部分。
- `qg2` (雇主性质): 控制不同部门间的收入差异。
- `qg6` (周工作小时): 控制工作强度的影响。
- `qa301` (户口状况): 户籍制度对就业和收入有重要影响。
- `provcd22` (省份代码): 用于构建省份固定效应，控制地区不随时间变化的异质性。
- `qv102` (父亲教育程度) 和 `qv202` (母亲教育程度): 控制家庭背景因素。
- 研究将使用 `employ` (当前工作状态) 筛选在业样本，并考虑使用 `rswt_natcs22n` (横截面权重) 进行加权估计以提高样本代表性。

识别策略：主要依赖OLS模型的控制变量法来缓解遗漏变量偏误。通过纳入省份固定效应控制地区层面的不可观测因素。将使用稳健标准误来处理潜在的异方差问题。教育的内生性是本研究的一个潜在挑战，未来可考虑工具变量法。





## 回归结果

### WLS Regression Results

```
=====
Dep. Variable:          log_emp_income    R-squared:                0.429
Model:                  WLS              Adj. R-squared:          0.350
Method:                 Least Squares     F-statistic:              341.6
Date:                   Thu, 29 May 2025  Prob (F-statistic):      2.09e-276
Time:                   13:14:08          Log-Likelihood:           -610.80
No. Observations:       421              AIC:                     1326.
Df Residuals:           369              BIC:                     1536.
Df Model:                51
Covariance Type:        HCL
=====
```

	coef	std err	t	P> t	[0.025	0.975]
Intercept	7.2221	1.190	6.071	0.000	4.883	9.561
C(qea0)[T.2.0]	0.3553	0.176	2.016	0.045	0.009	0.702
C(qea0)[T.3.0]	0.4250	0.337	1.261	0.208	-0.238	1.088
C(qea0)[T.4.0]	0.1027	0.274	0.375	0.708	-0.435	0.641
C(qea0)[T.5.0]	-1.0119	0.362	-2.794	0.005	-1.724	-0.300
C(qg2)[T.2.0]	-0.5665	0.393	-1.442	0.150	-1.339	0.206
C(qg2)[T.3.0]	0.0675	0.314	0.215	0.830	-0.550	0.685
C(qg2)[T.4.0]	-0.1602	0.299	-0.536	0.592	-0.748	0.427
C(qg2)[T.5.0]	-0.6179	0.419	-1.475	0.141	-1.442	0.206
C(qg2)[T.6.0]	-0.6357	0.350	-1.816	0.070	-1.324	0.053
C(qg2)[T.7.0]	-0.3810	0.365	-1.044	0.297	-1.099	0.337
C(qg2)[T.8.0]	-0.7827	0.405	-1.931	0.054	-1.580	0.014
C(qa301)[T.3.0]	-0.0077	0.186	-0.042	0.967	-0.373	0.357
C(qa301)[T.7.0]	0.1702	0.209	0.814	0.416	-0.241	0.581
C(provcd22)[T.12.0]	-0.8030	1.088	-0.738	0.461	-2.942	1.336
C(provcd22)[T.13.0]	-0.5336	0.905	-0.590	0.556	-2.313	1.246
C(provcd22)[T.14.0]	0.1889	0.922	0.205	0.838	-1.624	2.001
C(provcd22)[T.21.0]	0.0020	0.854	0.002	0.998	-1.677	1.681

C(provcd22) [T.22.0]	0.2698	0.884	0.305	0.760	-1.469	2.009
C(provcd22) [T.23.0]	-0.9532	1.061	-0.898	0.370	-3.040	1.134
C(provcd22) [T.31.0]	0.5277	0.876	0.603	0.547	-1.194	2.249
C(provcd22) [T.32.0]	0.1886	0.959	0.197	0.844	-1.698	2.075
C(provcd22) [T.33.0]	0.3749	0.906	0.414	0.679	-1.406	2.156
C(provcd22) [T.34.0]	0.3135	0.902	0.348	0.728	-1.459	2.086
C(provcd22) [T.35.0]	0.4056	0.891	0.455	0.649	-1.347	2.158
C(provcd22) [T.36.0]	0.2587	0.937	0.276	0.783	-1.583	2.101
C(provcd22) [T.37.0]	0.1996	0.914	0.218	0.827	-1.598	1.997
C(provcd22) [T.41.0]	-0.4198	0.919	-0.457	0.648	-2.227	1.387
C(provcd22) [T.42.0]	-0.6810	0.999	-0.682	0.496	-2.645	1.283
C(provcd22) [T.43.0]	0.1512	0.871	0.174	0.862	-1.561	1.864
C(provcd22) [T.44.0]	0.3713	0.872	0.426	0.670	-1.343	2.085
C(provcd22) [T.45.0]	-0.1491	0.903	-0.165	0.869	-1.925	1.626
C(provcd22) [T.50.0]	-0.7229	1.153	-0.627	0.531	-2.991	1.545
C(provcd22) [T.51.0]	-0.0931	0.916	-0.102	0.919	-1.894	1.708
C(provcd22) [T.52.0]	-0.1298	0.901	-0.144	0.886	-1.902	1.643
C(provcd22) [T.53.0]	0.4166	0.894	0.466	0.642	-1.341	2.175
C(provcd22) [T.54.0]	1.0799	0.882	1.224	0.222	-0.655	2.815
C(provcd22) [T.61.0]	-0.4541	0.939	-0.484	0.629	-2.300	1.392
C(provcd22) [T.62.0]	-0.6103	0.924	-0.660	0.510	-2.428	1.207
C(provcd22) [T.63.0]	-0.8140	1.247	-0.653	0.514	-3.265	1.637
C(provcd22) [T.64.0]	0.6853	0.899	0.762	0.446	-1.082	2.453
C(provcd22) [T.65.0]	-0.3076	1.006	-0.306	0.760	-2.285	1.670
cfps2022eduy_im	0.0850	0.021	3.960	0.000	0.043	0.127
age	0.0594	0.034	1.731	0.084	-0.008	0.127
age_sq	-0.0008	0.000	-1.998	0.046	-0.002	-1.27e-05
gender	0.5759	0.115	5.009	0.000	0.350	0.802
urban22	0.2812	0.135	2.083	0.038	0.016	0.547
party	0.1700	0.226	0.752	0.453	-0.275	0.615
qp201	0.1220	0.052	2.328	0.020	0.019	0.225
qg6	0.0086	0.005	1.836	0.067	-0.001	0.018
qv102	0.0090	0.012	0.771	0.441	-0.014	0.032
qv202	0.0104	0.003	3.311	0.001	0.004	0.017

```

=====
Omnibus:                192.693    Durbin-Watson:                2.114
Prob(Omnibus):           0.000    Jarque-Bera (JB):            1457.420
Skew:                    -1.793    Prob(JB):                     0.00
Kurtosis:                11.380    Cond. No.                     7.80e+04
=====

```

#### Notes:

[1] Standard Errors are heteroscedasticity robust (HC1)

[2] The condition number is large, 7.8e+04. This might indicate that there are strong multicollinearity or other numerical problems.

#### Notes:

1. Standard errors are robust to heteroskedasticity (HC1 type).

2. The model was estimated using weighted Least Squares (WLS) with 'rswt\_natcs22n' as weights.

3. Categorical variables (qea0, qg2, qa301, provcd22) are included as dummy variables (fixed effects for provcd22).

The reference category for each C() variable is its lowest valid numerical code after cleaning.

## 结果解读

好的，作为一名计量经济学家，我对这份研究进行解读。

### 1. 各主要自变量对因变量的经济含义解释（因变量为对数化的“过去12个月所有工作（主要工作+一般工作）的税后工资性收入” `emp_income`）：

- “受教育年限” (`cfps2022eduy_im`)：系数为0.0850 ( $p < 0.001$ )，意味着在其他条件不变时，受教育年限每增加一年，个人税后工资性收入平均提高约8.5%。这体现了教育作为人力资本投资的正向回报。
- “年龄” (`age`) 与 “年龄平方” (`age_sq`)：年龄的线性项系数为0.0594 ( $p = 0.084$ )，平方项系数为-0.0008 ( $p = 0.046$ )。这共同表明工资性收入随年龄（作为工作经验的代理）的增加而增加，但其增长速度随年龄递减，呈现倒U型关系，符合人力资本理论中经验积累的边际效用递减规律。
- “受访者性别” (`gender`)：系数为0.5759 ( $p < 0.001$ )，以女性为参照组（编码0为女，1为男），表明男性的税后工资性收入平均比女性高出约57.6%，显示存在显著的性别工资差距。
- “城乡分类” (`urban22`)：系数为0.2812 ( $p = 0.038$ )，以乡村为参照组（编码0为乡村，1为城镇），表明城镇居民的税后工资性收入平均比乡村居民高出约28.1%，反映了城乡二元结构下的收入差异。
- 其他显著控制变量如“当前婚姻状态” (`qea0`) 中，“在婚（有配偶）”者（相对未婚）收入高35.5%，“丧偶”者收入低约63.6% ( $1 - \exp(-1.0119)$ )；“健康状况” (`qp201`) 每差一级（数值越大越不健康），收入反而高12.2%；“每周工作时间” (`qg6`) 每增加一小时，收入提高0.86%；“母亲教育程度” (`qv202`) 每提高一单位，子女收入增加1.04%。

### 2. 研究发现总结

本研究基于扩展的明瑟方程，采用加权最小二乘法（WLS）并控制省份固定效应，探讨了中国个体税后工资性收入（对数化）的决定因素。模型整体解释力尚可（调整后 $R^2$ 为0.350），F统计量显著。

核心发现如下：

- 教育回报显著**：“受教育年限”对收入有显著的正向影响，每多接受一年教育，收入平均提高8.5%，证实了教育的人力资本价值。
- 经验效应呈倒U型**：“年龄”及其平方项的系数表明，收入随工作经验的积累先上升后下降，符合生命周期理论。
- 性别与城乡差异**：存在显著的性别工资溢价，男性收入远高于女性。同时，城镇居民收入显著高于乡村居民，反映了持续存在的结构性不平等。

此外，控制变量分析显示：“当前婚姻状态”对收入有重要影响，“在婚”状态与较高收入相关，而“丧偶”则与较低收入相关。“母亲教育程度”对子女收入有正向的代际影响。“每周工作时间”作为劳动投入的衡量，其增加能带来更高的收入。值得注意的是，“健康状况”变量显示健康状况越差收入越高的反常现象（系数0.1220， $p = 0.020$ ），这可能源于未观测到的共同因素（如从事高风险高收入行业）、变量测量方式或样本选择偏误，需谨慎解读并作进一步研究。模型提示存在潜在的多重共线性问题（条件数较大），提示未来研究可关注此点。研究已使用稳健标准误差应对异方差。