

研究性别和收入之间的关系。

摘要与引言

本研究基于中国家庭追踪调查（CFPS）2022年数据（样本量8181），采用OLS回归分析探讨个人特征对税后工资性收入（ln_income）的影响，旨在揭示收入不平等的驱动因素。研究意义在于识别性别歧视、人力资本和区域差异对经济的贡献，为政策干预（如促进教育公平和反歧视）提供依据。方法上，控制年龄、教育、城乡、工作状态等变量，发现男性收入显著高于女性49.23%（ $p<0.01$ ），教育年限每增加一年收入提升7.45%（ $p<0.01$ ），城镇居民收入高12.93%（ $p<0.01$ ），区域差异（如上海收入更高）及健康、婚姻状态影响显著。模型拟合良好（ $R^2=0.268$ ），凸显多重因素交互作用。（字数：198）

研究计划

本研究探讨性别对个人年收入的影响。核心研究问题：在控制个人特征后，男性相对于女性的年收入是否存在显著差异？采用OLS回归模型，因收入为连续变量且分布右偏（emp_income的偏度7.9），建议对收入取对数处理。关键变量：因变量为emp_income（过去12个月所有工作税后收入），因其全面覆盖收入来源；核心自变量为gender（性别，男性=1），预期男性收入更高，依据性别工资差距理论（如劳动力市场歧视）。控制变量包括age（年龄，缓解工作经验差异）、cfps2022eduy（教育年限，人力资本理论）、employ（工作状态，区分就业状态）、urban22（城乡分类，控制区域经济差异）、marriage_last（婚姻状态，关联家庭责任）、qp201（健康状况，影响劳动力参与）和provcd22（省份固定效应），以减轻遗漏变量偏误。识别策略：加入provcd22省份固定效应，控制不可观测的地区异质性。

[图片 plot_20250711_124823_deepseek.png 未生成]

回归结果

OLS Regression Results						
=====						
Dep. Variable:	ln_income	R-squared:	0.268			
Model:	OLS	Adj. R-squared:	0.264			
Method:	Least Squares	F-statistic:	67.81			
Date:	Fri, 11 Jul 2025	Prob (F-statistic):	0.00			
Time:	13:12:50	Log-Likelihood:	-10720.			
No. Observations:	8181	AIC:	2.153e+04			
Df Residuals:	8136	BIC:	2.185e+04			
Df Model:	44					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

const	9.3672	0.127	73.724	0.000	9.118	9.616
age	-0.0122	0.001	-10.721	0.000	-0.014	-0.010
gender	0.4923	0.021	23.846	0.000	0.452	0.533
urban22	0.1293	0.022	5.780	0.000	0.085	0.173
cfps2022eduy	0.0745	0.003	24.860	0.000	0.069	0.080
employ_1.0	0.5283	0.079	6.649	0.000	0.373	0.684
employ_3.0	-0.0772	0.091	-0.850	0.395	-0.255	0.101
marriage_last_2.0	0.3155	0.030	10.643	0.000	0.257	0.374
marriage_last_3.0	0.2717	0.166	1.634	0.102	-0.054	0.598

marriage_last_4.0	0.3074	0.064	4.826	0.000	0.183	0.432
marriage_last_5.0	0.0970	0.094	1.031	0.303	-0.088	0.281
qp201_2.0	0.0284	0.035	0.822	0.411	-0.039	0.096
qp201_3.0	0.0438	0.029	1.516	0.129	-0.013	0.100
qp201_4.0	0.0337	0.050	0.676	0.499	-0.064	0.131
qp201_5.0	-0.1664	0.044	-3.745	0.000	-0.253	-0.079
provcd22_12.0	-0.1117	0.124	-0.903	0.367	-0.354	0.131
provcd22_13.0	-0.4746	0.089	-5.360	0.000	-0.648	-0.301
provcd22_14.0	-0.4921	0.093	-5.309	0.000	-0.674	-0.310
provcd22_15.0	-0.1454	0.283	-0.513	0.608	-0.701	0.410
provcd22_21.0	-0.4115	0.087	-4.703	0.000	-0.583	-0.240
provcd22_22.0	-0.4781	0.110	-4.365	0.000	-0.693	-0.263
provcd22_23.0	-0.5213	0.101	-5.141	0.000	-0.720	-0.323
provcd22_31.0	0.2408	0.091	2.646	0.008	0.062	0.419
provcd22_32.0	0.1578	0.100	1.576	0.115	-0.039	0.354
provcd22_33.0	0.1155	0.098	1.173	0.241	-0.077	0.309
provcd22_34.0	-0.0125	0.105	-0.118	0.906	-0.219	0.194
provcd22_35.0	-0.2218	0.124	-1.796	0.073	-0.464	0.020
provcd22_36.0	-0.4373	0.108	-4.048	0.000	-0.649	-0.226
provcd22_37.0	-0.2803	0.090	-3.115	0.002	-0.457	-0.104
provcd22_41.0	-0.4269	0.085	-4.993	0.000	-0.594	-0.259
provcd22_42.0	-0.1924	0.115	-1.669	0.095	-0.418	0.034
provcd22_43.0	-0.1643	0.099	-1.654	0.098	-0.359	0.030
provcd22_44.0	-0.1050	0.086	-1.225	0.221	-0.273	0.063
provcd22_45.0	-0.3938	0.110	-3.579	0.000	-0.609	-0.178
provcd22_46.0	0.2716	0.376	0.722	0.470	-0.466	1.009
provcd22_50.0	-0.2031	0.148	-1.370	0.171	-0.494	0.088
provcd22_51.0	-0.4082	0.095	-4.278	0.000	-0.595	-0.221
provcd22_52.0	-0.3863	0.105	-3.672	0.000	-0.592	-0.180
provcd22_53.0	-0.4567	0.103	-4.445	0.000	-0.658	-0.255
provcd22_54.0	0.8736	0.411	2.126	0.034	0.068	1.679
provcd22_61.0	-0.3545	0.103	-3.450	0.001	-0.556	-0.153
provcd22_62.0	-0.4425	0.086	-5.123	0.000	-0.612	-0.273
provcd22_63.0	-0.1971	0.272	-0.725	0.469	-0.730	0.336
provcd22_64.0	-0.1296	0.262	-0.494	0.621	-0.644	0.384
provcd22_65.0	-0.3265	0.138	-2.363	0.018	-0.597	-0.056

```

=====
Omnibus:                4578.997    Durbin-Watson:                1.955
Prob(Omnibus):           0.000    Jarque-Bera (JB):            89225.011
Skew:                    -2.273    Prob(JB):                     0.00
Kurtosis:                18.527    Cond. No.                    2.12e+03
=====

```

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.12e+03. This might indicate that there are strong multicollinearity or other numerical problems.

结果解读

任务1：自变量对因变量的经济含义（≤300字）

本研究因变量为“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”的自然对数（ln_income），系数解释为百分比变化。各自变量的经济含义如下（引用中文标签）：

- 受访者性别 (gender)**：系数0.4923 ($p<0.01$)，表示男性 (gender=1) 相比女性 (gender=0)，收入平均高49.23%，体现显著性别工资差距，可能源于劳动力市场歧视或职业隔离（例如男性在薪酬更高行业）。
- 年龄 (age)**：系数-0.0122 ($p<0.01$)，年龄每增加一岁，收入下降1.22%，反映生命周期效应（如经验回报递减或健康衰退影响生产力）。
- 基于国家统计局资料的城乡分类 (urban22)**：系数0.1293 ($p<0.01$)，城镇 (urban22=1) 相比乡村 (urban22=0)，收入高12.93%，归因于城镇更高经济发展水平和就业机会。
- CFPS2022个人问卷受访者已完成的受教育年限 (cfps2022eduy)**：系数0.0745 ($p<0.01$)，教育年限每增加一年，收入提高7.45%，符合人力资本理论（教育提升技能和收入潜力）。
- 当前工作状态 (employ)**：在业状态 (employ_1.0) 系数0.5283 ($p<0.01$)，相比参考组（失业或非经济活动），收入高52.83%，凸显就业对收入的积极贡献。
- 加载变量：最近一次访问婚姻状态 (marriage_last)**：在婚状态 (marriage_last_2.0) 系数0.3155 ($p<0.01$)，相比未婚，收入高31.55%，可能与家庭稳定增强经济保障相关。
- 健康状况 (qp201)**：不健康状态 (qp201_5.0) 系数-0.1664 ($p<0.01$)，相比参考组（非常健康），收入低16.64%，显示健康损害劳动力参与和生产力。
- 2022年省国标码 (provcd22)**：如上海 (provcd22_31.0) 系数0.2408 ($p<0.01$)，相比参考省份收入更高，体现区域经济差异（如沿海地区高收入机会）。

任务2：研究发现总结（≤400字）

本研究通过OLS回归分析性别对个人年收入的影响，控制年龄、教育、工作状态等变量，样本量8181。关键发现如下：

首先，性别工资差距显著：男性收入比女性高49.23% ($p<0.01$)，控制其他因素后，这一差距仍突出，支持性别歧视理论。这可能源于劳动力市场中的结构性偏见（如职业隔离或薪酬不公），凸显政策需关注性别平等。

其次，人力资本和区域因素驱动收入：教育年限每增加一年，收入提升7.45% ($p<0.01$)，证实教育投资回报高；城乡差异明显，城镇居民收入高12.93% ($p<0.01$)，反映城市化红利。此外，省份固定效应显示区域不平等：例如上海收入显著高于参考省份（系数0.2408, $p<0.01$ ），而河北、山西等系数负值大（如河北-0.4746, $p<0.01$ ），表明东部沿海地区优势。

第三，个人特征影响复杂：年龄增加导致收入年降1.22% ($p<0.01$)，可能与健康衰退相关；健康不佳者收入低16.64% ($p<0.01$)，强调健康人力资本重要性。工作状态中，在业者收入高52.83% ($p<0.01$)，但失业状态影响不显著；在婚者收入高31.55% ($p<0.01$)，暗示婚姻稳定性提升经济安全。

模型整体拟合良好 ($R^2=0.268$)，但需注意多重共线性风险（条件数 $2.12e+03$ ）。结论：性别是收入差距核心因素，但教育、区域和健康等交互作用强化不平等，建议政策干预如促进教育公平、区域协调发展和反歧视措施，以缩小差距。

