

股权特征与企业市场价值

吉小浩

2024-05-04

目 录

1. 引言.....	3
1.1 研究问题.....	3
1.2 研究意义.....	3
1.3 研究思路.....	4
2. 数据准备	4
2.1 数据来源.....	4
2.2 样本选择.....	4
2.3 合并数据.....	4
2.4 清洗数据	15
2.5 数据描述性统计.....	23
3. 模型训练与评估.....	38
3.1 算法基本介绍	38
3.2 模型训练.....	38
3.3 模型评估	41
4. 模型调参	41
5. 模型再训练与评估	44
6. 模型解释与应用	45
6.1 股权特征重要性.....	45
7. 研究结论	48
7.1 股权特征典型事实.....	48

说明

为方便大家做数据分析报告，邀请刘瑞雪和吉小浩两同学基于《课程作业模板》各做一份数据分析报告。

欢迎大家阅读、学习和参考，并结合自己的选题完成《Python 数据分析》课程数据分析报告。

《课程作业模板》在学习通中，请下载使用。

欢迎学习借鉴，鼓励创新，禁止抄袭。

授课教师：程 军

2024 年 5 月 5 日

本分析报告的实验环境设置如下：

```
! python --version
! pip list | findstr "pandas numpy scikit-learn"

Python 3.12.3
numpy          1.26.4
pandas         2.2.2
scikit-learn   1.4.2

from sklearn.model_selection import train_test_split, cross_val_score, RandomizedSearchCV
from IPython.core.interactiveshell import InteractiveShell
from sklearn.metrics import (
    median_absolute_error,
    mean_absolute_error,
    mean_squared_error,
    r2_score,
)
import pandas as pd
import numpy as np
from pprint import pprint
```

```
from IPython.core.interactiveshell import InteractiveShell

InteractiveShell.ast_node_interactivity = "all"

import matplotlib_inline

matplotlib_inline.backend_inline.set_matplotlib_formats("svg")

import warnings
warnings.filterwarnings('ignore')
```

1. 引言

1.1 研究问题

本分析报告主要关注 A 股上市公司的股权特征的典型事实及其对企业市场价值的影响。

股权特征：

- (1)两权分离度的变化趋势是怎样的？
- (2)机构投资者持股特征是怎样的？
- (3)国企与非国企，制造业与非制造业的股权集中度是怎样的？

与市值关系：

- (4)哪一项股权特征更为重要？

1.2 研究意义

股权特征反映了企业内部控制权的分配和股东之间的利益关系，而企业市值则是市场对企业整体价值的评估。通过研究股权特征与企业市值的关系，可以深入了解企业治理结构对企业价值的影响，以及不同股权结构下企业市值变动的规律。

通过研究，可以发现不同股权特征对企业市值的影响程度和机制。有助于企业根据自身的实际情况，优化股权结构，提高公司治理效率，从而提升企业市值。例如，合理调整国有股、法人股和社会公众股的比例，可以平衡不同股东之间的利益，增强企业的市场竞争力。

通过研究股权特征与企业市值的关系，可以为投资者提供重要的决策参考，帮助他们更准确地评估企业的投资价值和风险。

股权特征和企业市值之间的关系研究，有助于为政策制定者提供有益的参考，帮助他们制定更加科学合理的政策，推动资本市场的健康稳定发展。

1.3 研究思路

本报告旨在探讨股权特征对企业市值的影响。分析实际控制人性质、股东持股比例、机构持股情况等股权特征与企业市值之间的关系。同时，本研究将采用梯度提升决策树（GBDT）和 Lasso 回归模型进行模型分析。

2. 数据准备

2.1 数据来源

本报告使用的数据主要是上市公司中的财务数据和股权特征两部分。数据来自国泰安数据库（CAMAR）。

2.2 样本选择

数据区间：以数据起始时间为起点，截止到 2023 年。其中，参考常规做法，对于数据我们进行了 ST、金融业的剔除。在数据清理期间，填补可填补的数据，删除缺失较多的数据。最终形成 327567 个样本。

其中，

被解释变量为企业市场价值，本报告选择现有文献中常用指标托宾 Q 值。

解释变量即股权特征，选取机构持股数量，机构持股比例，股股东直接持有上市公司股份的数量，公司前 3 位大股东持股比例之和，两权分离率等。

控制变量，选择企业规模，无形资产占比，资产负债率，第一大股东持股比例之和等。

2.3 合并数据

2.3.1 财务指标数据

```
#财务指标
df_fi=pd.read_excel("rawdata\财务指标\FI_T10.xlsx",skiprows=[1,2],header=0)
df_fi.head()
```

	Stk	Short	Acc	In	Ind	F100	F10080	F100	F100	F100	F100
	cd	Name	per	dc	nme	101B	1A	901A	902A	903A	904A
0	2	深万 科 A	199 1- 12- 31	K7 0	房 地 产 业	NaN	1.50516 7e+09	2.647 639	2.648 996	3.243 873	3.245 536
1	2	深万	199	K7	房	NaN	2.44652	2.587	2.821	3.162	3.447

Stk cd	Short Name	Acc per	In dc d	Ind nme	F100 101B	F10080 1A	F100 901A	F100 902A	F100 903A	F100 904A
	科 A	2- 12- 31	0	地 产 业		2e+09	538	200	024	563
2 2	深万 科 A	199 3- 12- 31	K7 0	房 地 产 业	NaN	3.54123 1e+09	1.657 757	1.659 930	1.879 755	1.882 219
3 2	深万 科 A	199 4- 12- 31	K7 0	房 地 产 业	7.966 339	2.79743 5e+09	1.045 697	1.058 723	1.073 715	1.087 090
4 2	深万 科 A	199 5- 12- 31	K7 0	房 地 产 业	6.027 127	2.93386 2e+09	0.907 087	0.914 923	0.881 404	0.889 018

```
df_fi["year"]=df_fi["Accper"].str.slice(0,4).astype("int64")

#新建'manufacturingindustry', 根据行业名称, 含有制造业, 则赋值为 1, 其余赋值为 0
df_fi['manufacturingindustry'] = df_fi['Indnme'].str.contains('制造业',
na=False).astype(int)

df_fi.columns
df_fi.drop(["Accper"],axis=1,inplace=True)

Index(['Stkcd', 'ShortName', 'Accper', 'Indcd', 'Indnme', 'F100101B',
      'F100801A', 'F100901A', 'F100902A', 'F100903A', 'F100904A', 'year',
      'manufacturingindustry'],
      dtype='object')
```

2.3.2 股权特征数据

```
#股权特征
df_cs=pd.read_excel("rawdata\公司治理\HLD_Contrshr.xlsx",skiprows=[1,2],
header=0)
df_cr=pd.read_excel("rawdata\公司治理\HLD_CR.xlsx",skiprows=[1,2],header
=0)
df_hi=pd.read_excel("rawdata\公司治理\INI_Holder_IncrOrDesr.xlsx",skipro
ws=[1,2],header=0)

df_cs.head()
```

```
df_cr.head()
df_hi.head()
```

	Stkcd	Reptdt	S070 1a	S070 2a	S070 1b	S070 3a	S070 4a	S0705a	S070 2b	Seper ation
0	2	2003-12-31	1	0	0	NaN	NaN	NaN	NaN	NaN
1	2	2003-12-31	2	华润 股份 有限 公司	中国 华润 总公 司	国有 法人 股,募 集法 人股	11.1 9	156151 498.0	1100	NaN
2	2	2004-12-31	1	0	0	NaN	NaN	NaN	NaN	NaN
3	2	2004-12-31	2	华润 股份 有限 公司	中国 华润 总公 司	国有 法人 股,募 集法 人股, 流通 A 股	12.8 9	293084 169.0	1100	NaN
4	2	2005-12-31	1	0	0	NaN	NaN	NaN	NaN	NaN
	Stkcd		Reptdt		Shrcr1		Shrcr2		Shrcr3	
0	2		2003-12-31		11.1868		15.3752		17.4711	
1	2		2004-12-31		12.8906		17.1627		21.0274	
2	2		2005-12-31		11.8094		15.5075		18.5765	
3	2		2006-12-31		14.5428		18.8346		21.4614	
4	2		2007-12-31		14.6345		16.9825		19.0465	
	Institutio nID	Symbol		EndD ate	Statu s	HolderNumbe r		Shares		Proporti on
0	101775	2		1999-12-31	1	1		337676 5.0		0.61897 9
1	101775	2		1999-12-31	3	1		689200. 0		0.12633 4
2	101775	2		2000-12-31	0	1		620000 0.0		0.98261 1
3	101775	2		2000-12-31	2	1		228736 1.0		0.36251 4
4	101775	2		2000-	3	12		210081		3.32949

Stkcd	Reptdt	S070 1a	S070 2a	S070 1b	S070 3a	S070 4a	S0705a	S070 2b	Seper ation
			12-31				59.0	3	

```
df_cs=df_cs.rename(columns={'Reptdt':'Accper'})
df_cr=df_cr.rename(columns={'Reptdt':'Accper'})
df_hi=df_hi.rename(columns={'EndDate':'Accper','Symbol':'Stkcd'})

df_cs.head()
df_cr.head()
df_hi.head()
```

	Stkcd	Accper	S070 1a	S070 2a	S070 1b	S070 3a	S070 4a	S0705a	S070 2b	Seper ation
0	2	2003-12-31	1	0	0	NaN	NaN	NaN	NaN	NaN
1	2	2003-12-31	2	华润股份有限公司	中国华润总公司	国有法人股,募集法人股	11.19	156151498.0	1100	NaN
2	2	2004-12-31	1	0	0	NaN	NaN	NaN	NaN	NaN
3	2	2004-12-31	2	华润股份有限公司	中国华润总公司	国有法人股,募集法人股,流通A股	12.89	293084169.0	1100	NaN
4	2	2005-12-31	1	0	0	NaN	NaN	NaN	NaN	NaN
		Stkcd		Accper		Shrcr1		Shrcr2		Shrcr3
0		2		2003-12-31		11.1868		15.3752		17.4711
1		2		2004-12-31		12.8906		17.1627		21.0274
2		2		2005-12-31		11.8094		15.5075		18.5765
3		2		2006-12-31		14.5428		18.8346		21.4614
4		2		2007-12-31		14.6345		16.9825		19.0465
	Institutio nID	Stkcd		Accp er	Statu s	HolderNumbe r		Shares		Proporti on

	Stkcd	Accper	S070 1a	S070 2a	S070 1b	S070 3a	S070 4a	S0705a	S070 2b	Seper ation
0	101775	2		1999 -12- 31	1	1		337676 5.0	0.61897 9	
1	101775	2		1999 -12- 31	3	1		689200. 0	0.12633 4	
2	101775	2		2000 -12- 31	0	1		620000 0.0	0.98261 1	
3	101775	2		2000 -12- 31	2	1		228736 1.0	0.36251 4	
4	101775	2		2000 -12- 31	3	12		210081 59.0	3.32949 3	

```
df_sh1=pd.merge(df_cs,df_cr,how='outer',on=['Stkcd','Accper'])
df_sh=pd.merge(df_hi,df_sh1,how='outer',on=['Stkcd','Accper'])

df_sh.head()
```

	Ins titu tio nID	S t k c d	A c c p e r	S t a t u s	Hol der Nu mbe r	Sha res	Pr op ort ion	S 0 7 1 a	S 0 7 2 a	S 0 7 1 b	S 0 7 3 a	S 0 7 4 a	S 0 7 5 a	S e p e r a t i o n	S h e r e h o l d e r s	S h e r e h o l d e r s	S h e r e h o l d e r s
0	1017750	2	1	1.0	337676	0.61897	N	N	N	N	N	N	N	NaN	N	N	N
1	1017750	2	1	3.0	689200	0.12633	N	N	N	N	N	N	N	NaN	N	N	N

		S	A	t	Hol			S	S	S	S	S	S	S		S	S	S
Ins	t	cc	a	der		Pr	0	0	0	0	0	0	0	0	Se	h	h	h
titu	k	p	t	Nu		op	7	7	7	7	7	7	7	7	pe	rc	rc	rc
tio	c	e	u	mbe	Sha	ort	0	0	0	0	0	0	0	0	rat	r	r	r
nID	d	r	s	r	res	ion	1	2	1	3	4	5	2	ion	1	2	3	
		3																
		1																
2	10	2	2	0.	1.0	62	0.9	N	N	N	N	N	N	Na	N	N	N	
17			0	0		00	82	a	a	a	a	a	a	N	a	a	a	
75.			0			00	61	N	N	N	N	N	N		N	N	N	
0			0-			0.0	1											
		1																
		2-																
		3																
		1																
3	10	2	2	2.	1.0	22	0.3	N	N	N	N	N	N	Na	N	N	N	
17			0	0		87	62	a	a	a	a	a	a	N	a	a	a	
75.			0			36	51	N	N	N	N	N	N		N	N	N	
0			0-			1.0	4											
		1																
		2-																
		3																
		1																
4	10	2	2	3.	12.0	21	3.3	N	N	N	N	N	N	Na	N	N	N	
17			0	0		00	29	a	a	a	a	a	a	N	a	a	a	
75.			0			81	49	N	N	N	N	N	N		N	N	N	
0			0-			59.	3											
		1				0												
		2-																
		3																
		1																

#新建'PropertyRightsNature'列，根据实际控制人性质，若为国有性质，赋值为1【1100为国企，2000为行政机关、事业单位，2100为中央机构，2120为地方机构】，其中当实际控制人若有多人，只要其中之一是国有企业，判断为1；非国有性质，则赋值0

```
df_sh['PropertyRightsNature'] = df_sh['S0702b'].str.contains('1100|2000|2100|2120', case=False, na=False).astype(int)
```

```
df_sh["accper"]=df_sh["Accper"].str.slice(0,4).astype("int64")
```

```
df_sh.columns
```

```
df_sh.drop(["Accper"],axis=1,inplace=True)
```

```
Index(['InstitutionID', 'Stkcd', 'Accper', 'Status', 'HolderNumber', 'Shares',
```

```
      'Proportion', 'S0701a', 'S0702a', 'S0701b', 'S0703a', 'S0704a',
```

```
'S0705a', 'S0702b', 'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3',
'PropertyRightsNature', 'accper'],
dtype='object')
```

2.3.3 相关控制变量数据

```
df_cas=pd.read_excel('rawdata\控制变量\FS_Combas.xlsx',skiprows=[1,2],header=0)
df_cis=pd.read_excel('rawdata\控制变量\FS_Comins.xlsx',skiprows=[1,2],header=0)
df_ena=pd.read_excel('rawdata\控制变量\EN_EquityNatureAll.xlsx',skiprows=[1,2],header=0)
df_ipo=pd.read_excel('rawdata\控制变量\IPO_IPOCG.xlsx',skiprows=[1,2],header=0)
df_stk=pd.read_excel('rawdata\控制变量\STK_LISTEDCOINFOANL.xlsx',skiprows=[1,2],header=0)

df_cas.head()
df_cis.head()
df_ena.head()
df_ipo.head()
df_stk.head()
```

	Stkcd	Short Name	Accper	Typr	A00121 2000	A00121 8000	A00100 0000	A002000 000	A0032 00000
0	2	深万科 A	199 1- 12- 31	A	7.19523 4e+07	291259 .41	5.68494 1e+08	3.778345 e+08	9.2595 19e+04
1	2	深万科 A	199 2- 12- 31	A	9.41470 5e+07	783096 36.35	9.45501 6e+08	7.082929 e+08	1.7818 63e+06
2	2	深万科 A	199 3- 12- 31	A	1.65393 1e+08	279730 3.07	2.13615 9e+09	1.205853 e+09	1.8878 63e+06
3	2	深万科 A	199 4- 12- 31	A	2.99539 4e+08	329157 49.20	2.67518 8e+09	1.527296 e+09	5.8520 32e+07
4	2	深万科 A	199 5- 12- 31	A	2.47093 1e+08	277021 30.26	3.23437 7e+09	1.872121 e+09	1.8276 42e+08

	Stkcd	ShortName	Accper	Typrep	B001101000	B002000000
0	2	深万科 A	1991-12-31	A	NaN	NaN
1	2	深万科 A	1992-12-31	A	6.613562e+08	NaN
2	2	深万科 A	1993-12-31	A	1.084045e+09	NaN
3	2	深万科 A	1994-12-31	A	1.227544e+09	1.829864e+08
4	2	深万科 A	1995-12-31	A	1.503755e+09	1.668964e+08

	Symbol	ShortName	EndDate	LargestHolderRate	TopTenHolderRate
0	2	万科 A	2003-12-31	11.19	21.45
1	2	万科 A	2004-12-31	12.89	28.17
2	2	G 万科 A	2005-12-31	11.81	24.91
3	2	万科 A	2006-12-31	14.54	27.09
4	2	万科 A	2007-12-31	14.63	22.71

	Symbol	IPODate	DirectorNumber	IndDirectorNum	IsPartTime
0	1201	2021-04-14	5.0	NaN	1
1	1202	2021-04-19	10.0	6.0	1
2	1203	2021-04-20	3.0	NaN	0
3	1205	2021-04-28	8.0	NaN	0
4	1206	2021-05-06	10.0	3.0	1

	Symbol	ShortName	EndDate	EstablishDate
0	2	深万科 A	2000-12-31	1988-11-01
1	2	深万科 A	2001-12-31	1988-11-01
2	2	万科 A	2002-12-31	1988-11-01
3	2	万科 A	2003-12-31	1988-11-01
4	2	万科 A	2004-12-31	1988-11-01

合并表格

```
df_fs=pd.merge(df_cas,df_cis,how="outer",on=["Stkcd","Accper"])
df_fs["year"]=df_fs["Accper"].str.slice(0,4).astype("int64")

df_1=pd.merge(df_ena,df_stk,how="outer",on=["Symbol","EndDate"])
df_1=df_1.rename(columns={'EndDate':'Accper','Symbol':'Stkcd'})
df_1["accper"]=df_1["Accper"].str.slice(0,4).astype("int64")

df_ipo=df_ipo.rename(columns={'Symbol':'Stkcd'})
df_ipo["accper"]=df_ipo["IPODate"].str.slice(0,4).astype("int64")

df_fs.columns
df_fs.drop(["ShortName_x","Typrep_x","ShortName_y","Typrep_y"],axis=1,inplace=True)

Index(['Stkcd', 'ShortName_x', 'Accper', 'Typrep_x', 'A001212000',
       'A001218000', 'A001000000', 'A002000000', 'A003200000', 'ShortName_y',
       'Typrep_y', 'B001101000', 'B002000000', 'year'],
      dtype='object')

df_1.columns
df_1.drop(["ShortName_x","ShortName_y"],axis=1,inplace=True)

Index(['Stkcd', 'ShortName_x', 'Accper', 'LargestHolderRate',
       'TopTenHoldersRate', 'ShortName_y', 'EstablishDate', 'accper'],
      dtype='object')

df_ipo.columns

Index(['Stkcd', 'IPODate', 'DirectorNumber', 'IndDirectorNum', 'IsPartTime',
       'accper'],
      dtype='object')

df_kz=pd.merge(df_fs,df_1,how="left",left_on=["Stkcd","year"],right_on=["Stkcd","accper"])
df_kz=df_kz.merge(df_ipo,how="left",left_on=["Stkcd","year"],right_on=["Stkcd","accper"])

df_kz.head()
```

		A	A	A	A	A	B	B					E		Di	In		
	A	0	0	0	0	0	0	0		A	To	st	a	I	re	d	Is	a
	c	1	1	1	2	3	1	2		c	La	pT	a	c	ct	Di	P	c
S t k c d	c	2	2	0	0	2	1	0		c	rg	en	bl	c	P	or	a	c
	p	1	1	0	0	0	0	0	y	p	est	Ho	is	p	O	N	rt	p
	e	2	8	0	0	0	1	0	e	e	Ho	lde	h	e	D	u	T	e
	r	0	0	0	0	0	0	0	a	_	lde	rs	D	_	a	N	i	r
	_x	0	0	0	0	0	0	0	r	y	ate	ate	e	x	e	r	m	y
0	2	1	7.	2	5.	3.	9.	N	N	1	N	Na	Na	N	N	Na	Na	N
		9	1	9	6	7	2	a	a	9	a	N	N	a	a	N	N	a
		9	9	1	8	7	5	N	N	9	N			N	N		N	N
		1	5	2	4	8	9			1								
	-	2	5	9	3	5												
	1	3	9.	4	4	1												
	2	4	4	1	5	9												
	-	e	1	e	e	e												
	3	+		+	+	+												
	1	0		0	0	0												
		7		8	8	4												
1	2	1	9.	7	9.	7.	1.	6.	N	1	N	Na	Na	N	N	Na	Na	N
		9	4	8	4	0	7	6	a	9	a	N	N	a	a	N	N	a
		9	1	3	5	8	8	1	N	9	N			N	N		N	N
	2	4	0	5	2	1	3			2								
	-	7	9	0	9	8	5											
	1	0	6	1	2	6	6											
	2	5	3	6	9	3	2											
	-	e	6.	e	e	e	e											
	3	+	3	+	+	+	+											
	1	0	5	0	0	0	0											
		7		8	8	6	8											
2	2	1	1.	2	2.	1.	1.	1.	N	1	N	Na	Na	N	N	Na	Na	N
		9	6	7	1	2	8	0	a	9	a	N	N	a	a	N	N	a
		9	5	9	3	0	8	8	N	9	N			N	N		N	N
	3	3	7	6	5	7	4			3								
	-	9	3	1	8	8	0											
	1	3	0	5	5	6	4											
	2	1	3.	9	3	3	5											
	-	e	0	e	e	e	e											
	3	+	7	+	+	+	+											
	1	0		0	0	0	0											
		8		9	9	6	9											
3	2	1	2.	3	2.	1.	5.	1.	1.	1	N	Na	Na	N	N	Na	Na	N

	A	A	A	A	A	B	B					E		Di	In		
	0	0	0	0	0	0	0		A	To	st	a		re	d	Is	a
A	0	0	0	0	0	0	0		c	La	pT	a	c	I	ct	Di	P
c	1	1	1	2	3	1	2		c	rg	en	bl	c	P	or	re	a
c	2	2	0	0	2	1	0		p	est	Ho	is	p	O	N	ct	rt
S	1	1	0	0	0	0	0		y	Ho	lde	h	e	D	u	or	T
t	2	8	0	0	0	1	0		e	lde	rs	D	r	a	m	N	i
k	0	0	0	0	0	0	0		a	rR	Ra	at	_	t	be	u	m
c	0	0	0	0	0	0	0		r	y	ate	e	x	e	r	m	e
d	0	0	0	0	0	0	0		r	y	ate	e	x	e	r	m	e
	9	9	2	6	5	8	2	8	9	a	N	N	a	a	a	N	a
	9	9	9	7	2	5	2	2	9	N		N	N	N		N	N
	4	5	1	5	7	2	7	9	4								
	-	3	5	1	2	0	5	8									
	1	9	7	8	9	3	4	6									
	2	4	4	8	6	2	4	4									
	-	e	9.	e	e	e	e	e									
	3	+	2	+	+	+	+	+									
	1	0	0	0	0	0	0	0									
	8		9	9	7	9	8										
4	2	1	2.	2	3.	1.	1.	1.	1	N	Na	Na	N	N	N	Na	Na
	9	4	7	2	8	8	5	6	9	a	N	N	a	a	a	N	N
	9	7	7	3	7	2	0	6	9	N			N	N	N		N
	5	0	0	4	2	7	3	8	5								
	-	9	2	3	1	6	7	9									
	1	3	1	7	2	4	5	6									
	2	1	3	7	1	2	5	4									
	-	e	0.	e	e	e	e	e									
	3	+	2	+	+	+	+	+									
	1	0	6	0	0	0	0	0									
	8		9	9	8	9	8										

```
df_kz.columns
```

```
df_kz.drop(['Accper_x','Accper_y','accper_x','accper_y'],axis=1,inplace=True)
```

```
Index(['Stkcd', 'Accper_x', 'A001212000', 'A001218000', 'A001000000',
       'A002000000', 'A003200000', 'B001101000', 'B002000000', 'year',
       'Accper_y', 'LargestHolderRate', 'TopTenHoldersRate', 'Establish
Date',
       'accper_x', 'IPODate', 'DirectorNumber', 'IndDirectorNum', 'IsPa
rtTime',
       'accper_y'],
      dtype='object')
```

2.3.4 合并数据

```
#合并数据
df=pd.merge(df_fi,df_sh,how='left',left_on=["Stkcd","year"],right_on=["Stkcd","accper"])
df=df.merge(df_kz,how='left',left_on=["Stkcd","year"],right_on=["Stkcd","year"])
```

2.4 清洗数据

```
df.columns
Index(['Stkcd', 'ShortName', 'Indcd', 'Indnme', 'F100101B', 'F100801A',
      'F100901A', 'F100902A', 'F100903A', 'F100904A', 'year',
      'manufacturingindustry', 'InstitutionID', 'Status', 'HolderNumbe
r',
      'Shares', 'Proportion', 'S0701a', 'S0702a', 'S0701b', 'S0703a',
      'S0704a', 'S0705a', 'S0702b', 'Seperation', 'Shrcr1', 'Shrcr2',
      'Shrcr3', 'PropertyRightsNature', 'accper', 'A001212000', 'A0012
18000',
      'A001000000', 'A002000000', 'A003200000', 'B001101000', 'B002000
000',
      'LargestHolderRate', 'TopTenHoldersRate', 'EstablishDate', 'IPOD
ate',
      'DirectorNumber', 'IndDirectorNum', 'IsPartTime'],
      dtype='object')
df.head()
```

										A	B	B			E	Di	In	Is
										0	0	0	La	To	st	re	d	Is
										0	0	0	rg	pT	a	I	Di	P
										3	1	2	est	en	bl	P	or	re
S	h			I	0	0	0	0	0	0	0	0	Ho	Ho	is	O	N	or
t	o	N	n	d	1	8	9	9	9	9	0	1	ld	lde	h	D	u	or
k		a	d	n	0	0	0	0	0	0	0	0	er	rs	D	a	m	T
c		m	c	m	1	1	1	2	3	4	0	0	Ra	Ra	at	t	be	N
d	e	d	e	B	A	A	A	A	A	A	0	0	te	te	e	e	r	m
0	2	深	K	房	N	1.	2	2	3	3	9.	N	Na	Na	N	N	N	N
		万	7	地	a	5	2	a	N	N	a	a	a	a
		科	0	产	N	0	6	6	2	2	5	N			N	N	N	N
		A		业		5	4	4	4	4	9							
						1	7	8	3	5	5							
						6	6	9	8	5	1							
						7	3	9	7	3	9							
						e	9	6	3	6	e							
						+					+							
						0					0							

										A	B	B			E	Di	In	Is		
										0	0	0	La	To	st	I	re	d	Is	
										0	0	0	rg	pT	a	P	ct	Di	P	
										3	1	2	est	en	bl	O	or	re	a	
										2	1	0	Ho	Ho	is	D	N	ct	rt	
										0	1	0	ld	lde	h	a	u	or	T	
										0	0	0	er	rs	D	t	m	N	i	
										0	0	0	Ra	Ra	ate	e	be	u	m	
										0	0	0	te	te	e	e	r	m	e	
										9										
1	2	深	K	房	N	2.	2	2	3	3	.	1.	6.	N	Na	Na	N	N	N	
		万	7	地	a	4	7	6	a	N	N	a	a	a	
		科	0	产	N	4	5	8	1	4	.	8	1	N			N	N	N	
		A				6	8	2	6	4		1	3							
						5	7	1	2	7		8	5							
						2	5	2	0	5		6	6							
						2	3	0	2	6		3	2							
						e	8	0	4	3		e	e							
						+						+	+							
						0						0	0							
						9						6	8							
2	2	深	K	房	N	3.	1	1	1	1	.	1.	1.	N	Na	Na	N	N	N	
		万	7	地	a	5	8	0	a	N	N	a	a	a	
		科	0	产	N	4	6	6	8	8	.	8	8	N			N	N	N	
		A				1	5	5	7	8		7	4							
						2	7	9	9	2		8	0							
						3	7	9	7	2		6	4							
						1	5	3	5	1		3	5							
						e	7	0	5	9		e	e							
						+						+	+							
						0						0	0							
						9						6	9							
3	2	深	K	房	7	2.	1	1	1	1	.	5.	1.	Na	Na	N	N	N	N	
		万	7	地	.	7	8	2	N	N	a	a	a	a	
		科	0	产	9	9	0	0	0	0	.	5	2			N	N	N	N	
		A			6	7	4	5	7	8		2	7							
					6	4	5	8	3	7		0	5							
					3	3	6	7	7	0		3	4							
					3	5	9	2	1	9		2	4							
					9	e	7	3	5	0		e	e							
						+						+	+							
						0						0	0							
						9						7	9							

										A	B	B							
										0	0	0							
										0	0	0	La	To	E		Di	In	Is
										3	1	2	rg	pT	a	I	re	d	P
										2	1	0	est	en	bl	P	or	re	a
										0	0	0	Ho	Ho	is	O	N	ct	rt
										0	1	0	ld	lde	h	D	u	or	T
										0	0	0	er	rs	D	a	m	N	i
										0	0	0	Ra	Ra	at	t	be	u	m
										0	0	0	te	te	e	e	r	m	e
4	2	深	K	房	6	2.	0	0	0	0	1.	1.	1.	Na	Na	N	N	N	N
		万	7	地	.	9	8	5	6	N	N	a	a	a	a
		科	0	产	0	3	9	9	8	8	2	0	6			N	N	N	N
		A		业	2	3	0	1	8	8	7	3	8						
					7	8	7	4	1	9	6	7	9						
					1	6	0	9	4	0	4	5	6						
					2	2	8	2	0	1	2	5	4						
					7	e	7	3	4	8	e	e	e						
					+						+	+	+						
					0						0	0	0						
					9						8	9	8						

#清洗前样本量

`df.shape`

`(362440, 44)`

#删除 ShortName 中含有 s,st,*st 的样本

`df = df[~df['ShortName'].str.contains(r's|st|*st', case=False)]`

#删除 s,st,*st 的样本后

`df.shape`

`(354324, 44)`

`df.isnull().sum()`

```

Stkcd          0
ShortName      0
Indcd          0
Indnme         0
F100101B      34353
F100801A       0
F100901A       0
F100902A       0
F100903A       0
F100904A       0
year           0
manufacturingindustry  0

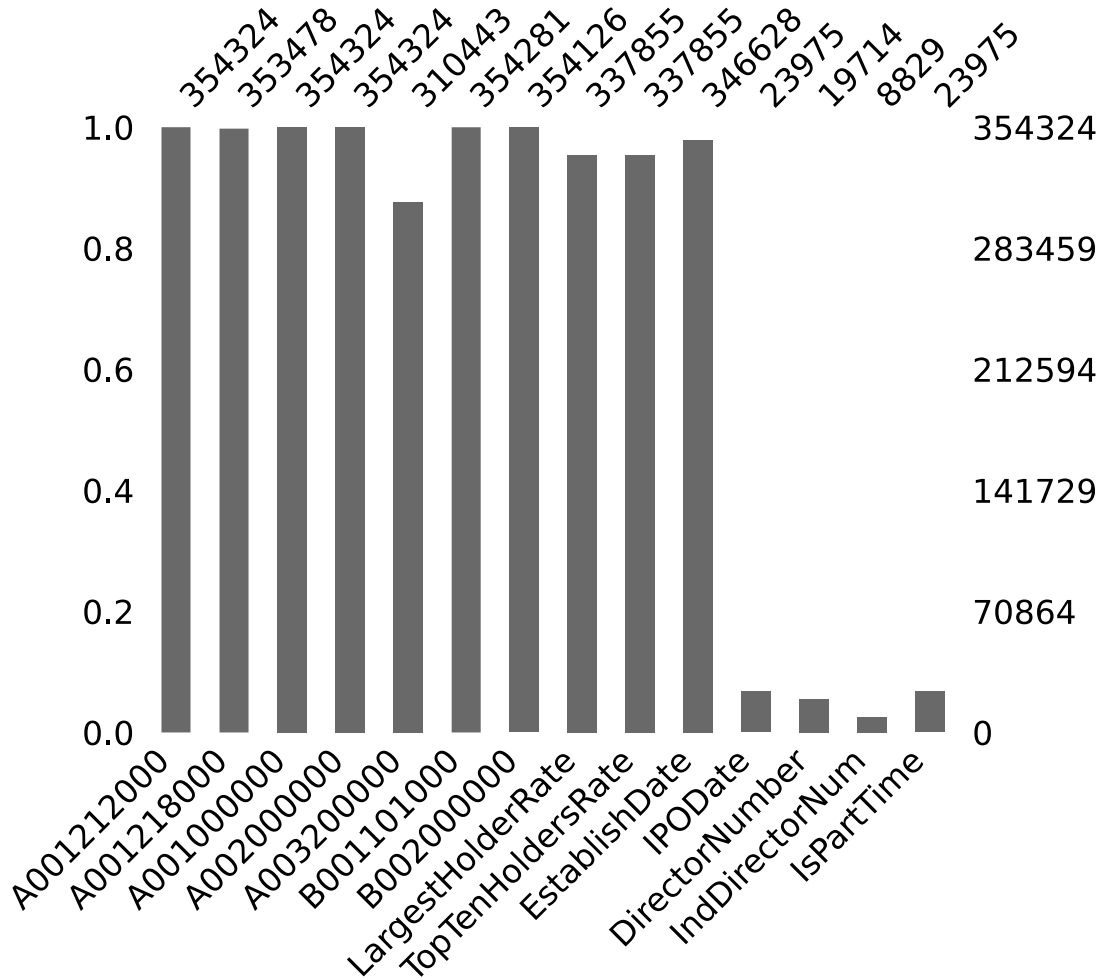
```

InstitutionID	4587
Status	4587
HolderNumber	4587
Shares	4587
Proportion	4587
S0701a	17443
S0702a	17443
S0701b	17443
S0703a	23660
S0704a	23694
S0705a	23713
S0702b	23737
Seperation	31493
Shrcr1	6974
Shrcr2	6974
Shrcr3	6974
PropertyRightsNature	4402
accper	4402
A001212000	0
A001218000	846
A001000000	0
A002000000	0
A003200000	43881
B001101000	43
B002000000	198
LargestHolderRate	16469
TopTenHoldersRate	16469
EstablishDate	7696
IPODate	330349
DirectorNumber	334610
IndDirectorNum	345495
IsPartTime	330349
dtype:	int64

df.columns

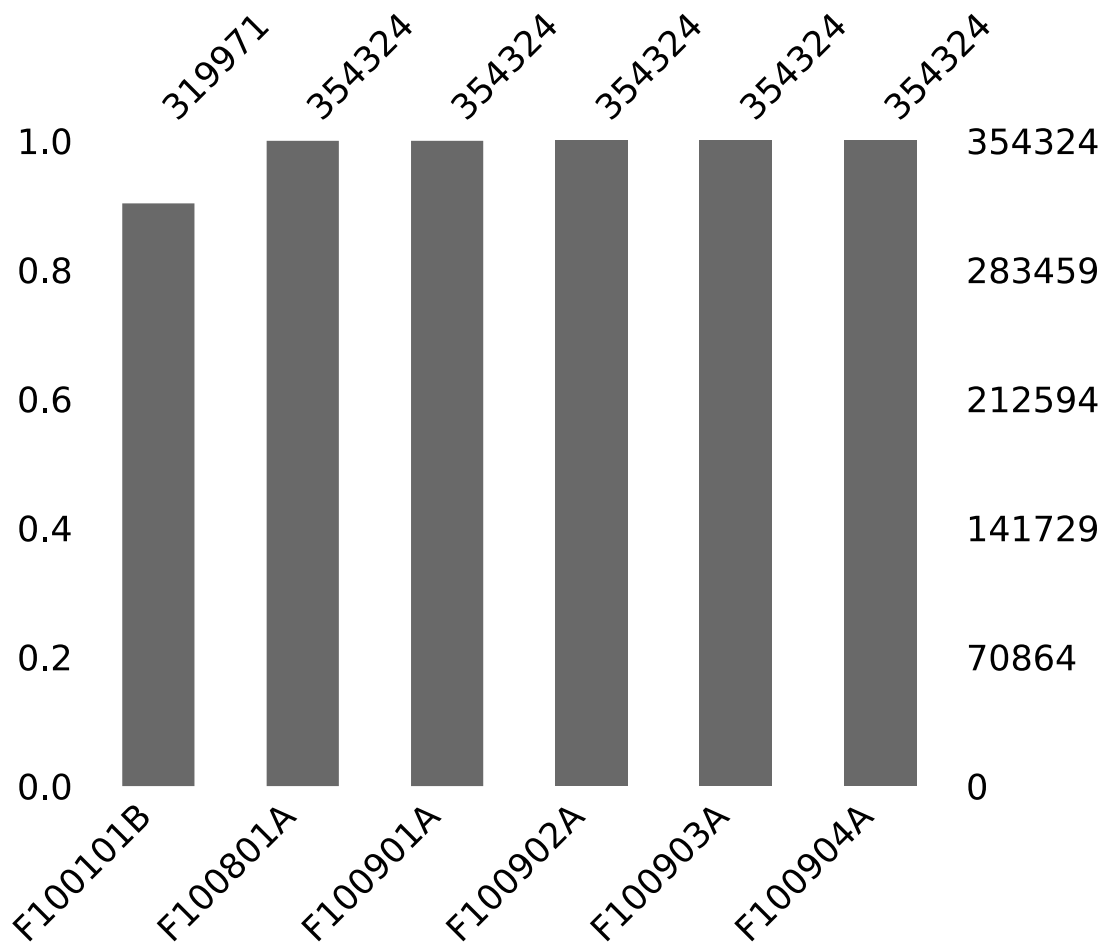
```
Index(['Stkcd', 'ShortName', 'Indcd', 'Indnme', 'F100101B', 'F100801A',
      'F100901A', 'F100902A', 'F100903A', 'F100904A', 'year',
      'manufacturingindustry', 'InstitutionID', 'Status', 'HolderNumber',
      'Shares', 'Proportion', 'S0701a', 'S0702a', 'S0701b', 'S0703a',
      'S0704a', 'S0705a', 'S0702b', 'Seperation', 'Shrcr1', 'Shrcr2',
      'Shrcr3', 'PropertyRightsNature', 'accper', 'A001212000', 'A001218000',
      'A001000000', 'A002000000', 'A003200000', 'B001101000', 'B002000000',
      'LargestHolderRate', 'TopTenHoldersRate', 'EstablishDate', 'IPODate',
      'DirectorNumber', 'IndDirectorNum', 'IsPartTime'],
      dtype='object')
```

```
import missingno as msno
msno.bar(df[['A001212000', 'A001218000',
              'A001000000', 'A002000000', 'A003200000', 'B001101000', 'B002000
000',
              'LargestHolderRate', 'TopTenHoldersRate', 'EstablishDate', 'IPOD
ate',
              'DirectorNumber', 'IndDirectorNum', 'IsPartTime']],figsize=(8,6),
fontSize=18)
```



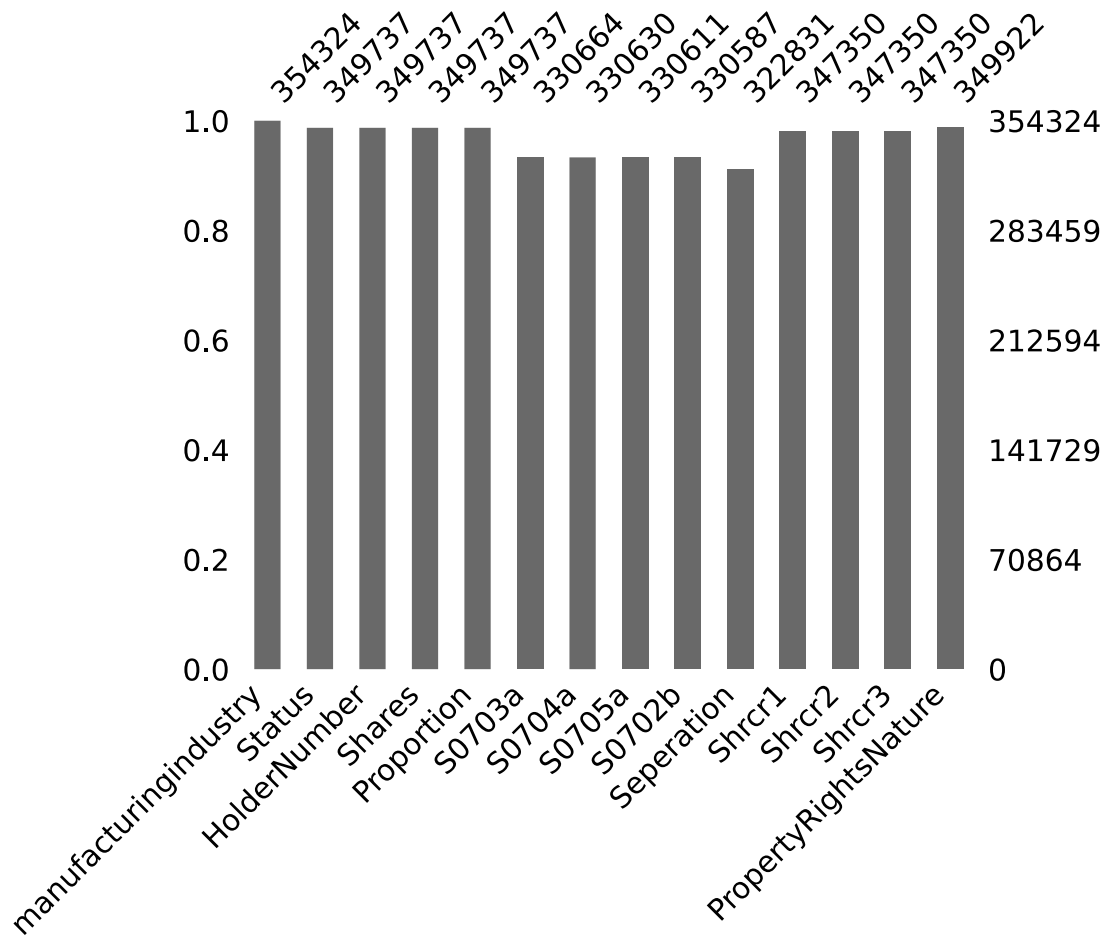
#市值指标缺失值可视化

```
import missingno as msno
msno.bar(df=df[['F100101B', 'F100801A',
                'F100901A', 'F100902A', 'F100903A', 'F100904A']], figsize=(8,6), fontsize=18)
```



#股权特征缺失值可视化

```
import missingno as msno
msno.bar(df=df[['manufacturingindustry','Status', 'HolderNumber', 'Shares', 'Proportion', 'S0703a', 'S0704a', 'S0705a', 'S0702b', 'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNature']],figsize=(8,6),fontsize=18)
```



```
# F100101B [市盈率 (PE) 1] 缺失过多, 删除
df.drop(["F100101B"],axis=1,inplace=True)
```

```
#填补股权特征中实际股权性质, 两权分离度, 直接控股股东持股比例, 机构持股比例缺失值
```

```
df["PropertyRightsNature"]=df["PropertyRightsNature"].fillna(0)
df["Seperation"]=df["Seperation"].fillna(0)
df["S0704a"]=df["S0704a"].fillna(0)
df["Proportion"]=df["Proportion"].fillna(0)
```

```
#其余缺失, 直接删除
df=df.dropna()
```

```
df.isnull().sum()
```

```
Stkcd          0
ShortName      0
Indcd          0
Indnme        0
F100801A      0
F100901A      0
F100902A      0
```

```

F100903A      0
F100904A      0
year          0
manufacturingindustry  0
InstitutionID  0
Status        0
HolderNumber  0
Shares        0
Proportion    0
S0701a        0
S0702a        0
S0701b        0
S0703a        0
S0704a        0
S0705a        0
S0702b        0
Seperation    0
Shrcr1        0
Shrcr2        0
Shrcr3        0
PropertyRightsNature  0
accper        0
A001212000    0
A001218000    0
A001000000    0
A002000000    0
B001101000    0
B002000000    0
LargestHolderRate  0
TopTenHoldersRate  0
EstablishDate  0
dtype: int64

```

计算控制变量

```
df.columns
```

```

Index(['Stkcd', 'ShortName', 'Indcd', 'Indnme', 'F100801A', 'F100901A',
      'F100902A', 'F100903A', 'F100904A', 'year', 'manufacturingindustry',
      'InstitutionID', 'Status', 'HolderNumber', 'Shares', 'Proportion',
      'S0701a', 'S0702a', 'S0701b', 'S0703a', 'S0704a', 'S0705a', 'S0702b',
      'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNature',
      'accper', 'A001212000', 'A001218000', 'A001000000', 'A002000000',
      'B001101000', 'B002000000', 'LargestHolderRate', 'TopTenHoldersRate'],
      dtype='object')

```

```

        'EstablishDate'],
        dtype='object')

df['lnAssets']=np.log(df['A001000000'])
df['Lev']=df['A002000000']/df['A001000000']
df['PPE_TA']=df['A001212000']/df['A001000000']
df['Intass']=df['A001218000']/df['A001000000']
df['OwnCon1']=df['LargestHolderRate']
df['OwnCon2_10']=df['TopTenHoldersRate']-df['LargestHolderRate']
df['SOE']=df['PropertyRightsNature']
df['ROA']=df['B002000000']/df['A001000000']

df["Establish_year"]=df["EstablishDate"].str.slice(0,4).astype("int64")
df['Age']=df['year']-df['Establish_year']

#.shift(1) 用于将数据沿轴（默认是行轴，即 axis=0）移动指定的数量。在这里，1
表示向下移动一行
df['Last_Year_Revenue'] = df['B001101000'].shift(1)
df['SalesGrowth'] = (df['B001101000'] / df['Last_Year_Revenue'] - 1).fillna(0)

```

2.5 数据描述性统计

```

df.columns

Index(['Stkcd', 'ShortName', 'Indcd', 'Indnme', 'F100801A', 'F100901A',
       'F100902A', 'F100903A', 'F100904A', 'year', 'manufacturingindust
ry',
       'InstitutionID', 'Status', 'HolderNumber', 'Shares', 'Proportion
',
       'S0701a', 'S0702a', 'S0701b', 'S0703a', 'S0704a', 'S0705a', 'S07
02b',
       'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNatur
e',
       'accper', 'A001212000', 'A001218000', 'A001000000', 'A002000000',
       'B001101000', 'B002000000', 'LargestHolderRate', 'TopTenHoldersR
ate',
       'EstablishDate', 'lnAssets', 'Lev', 'PPE_TA', 'Intass', 'OwnCon1
',
       'OwnCon2_10', 'SOE', 'ROA', 'Establish_year', 'Age',
       'Last_Year_Revenue', 'SalesGrowth'],
      dtype='object')

```

2.5.1 全描述性统计

```

#最终样本量
df.shape

(327567, 50)

df.describe()

```

[illegible]

</																					

2.5.2 描述性统计分析可视化

2.5.2.1 两权分离度可视化探究

```
import matplotlib.pyplot as plt

period1 = (2003, 2013)
period2 = (2013, 2023)

mean_separation1 = df[(df['year'] >= period1[0]) & (df['year'] <= period1[1])].groupby('year')['Seperation'].mean()
mean_separation2 = df[(df['year'] >= period2[0]) & (df['year'] <= period2[1])].groupby('year')['Seperation'].mean()

# 设置中文显示
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

#2003-2013
plt.figure(figsize=(10, 6))
plt.subplot(2, 1, 1)
plt.plot(mean_separation1.index, mean_separation1.values, marker='o', label='2003-2013')
plt.xlabel('年份')
plt.ylabel('两权分离度年度均值')
plt.title('2003-2013 两权分离度年度均值趋势')
plt.legend()

#2013-2023
plt.subplot(2, 1, 2)
plt.plot(mean_separation2.index, mean_separation2.values, marker='o', label='2013-2023')
plt.xlabel('年份')
plt.ylabel('两权分离度年度均值')
plt.title('2013-2023 年两权分离度年度均值趋势')
plt.legend()

plt.tight_layout() # 调整子图参数，使之填充整个图像区域
plt.show()

Text(0.5, 0, '年份')

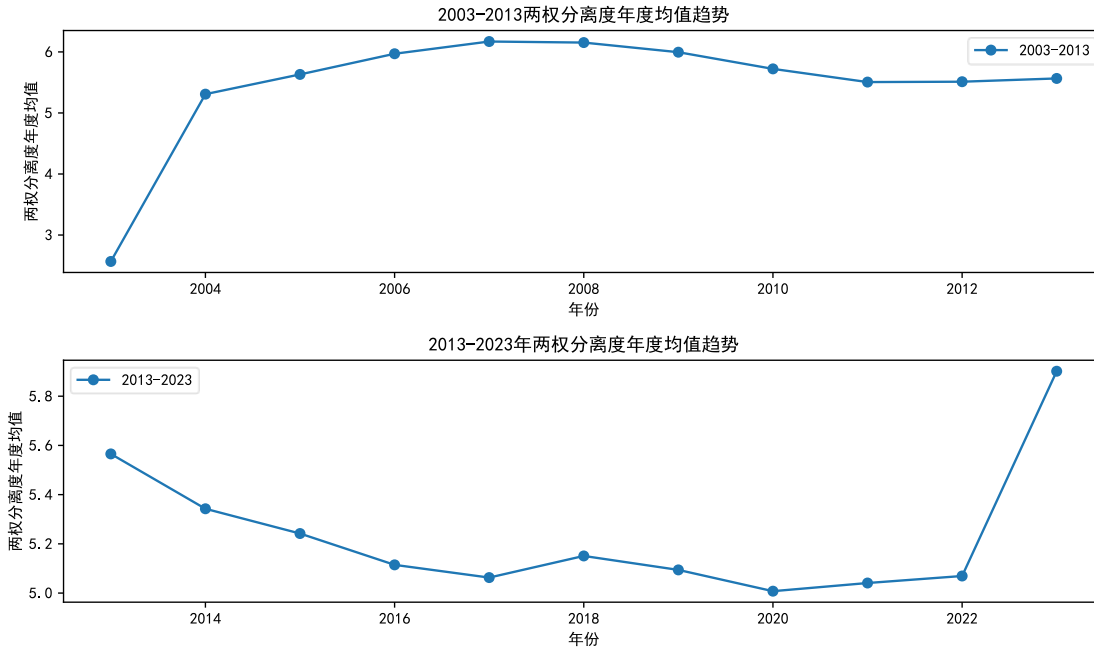
Text(0, 0.5, '两权分离度年度均值')

Text(0.5, 1.0, '2003-2013 两权分离度年度均值趋势')

Text(0.5, 0, '年份')

Text(0, 0.5, '两权分离度年度均值')
```

Text(0.5, 1.0, '2013-2023 年两权分离度年度均值趋势')



```
import matplotlib.pyplot as plt

means = df.groupby(['year', 'PropertyRightsNature'])['Seperation'].mean()
().reset_index()
yes_means = means[means['PropertyRightsNature'] == 1]
no_means = means[means['PropertyRightsNature'] == 0]

# 绘制国企的数据
plt.plot(yes_means['year'], yes_means['Seperation'], label='国企', marker='p', color='c')

# 绘制非国企的数据
plt.plot(no_means['year'], no_means['Seperation'], label='非国企', marker='d', color='g')

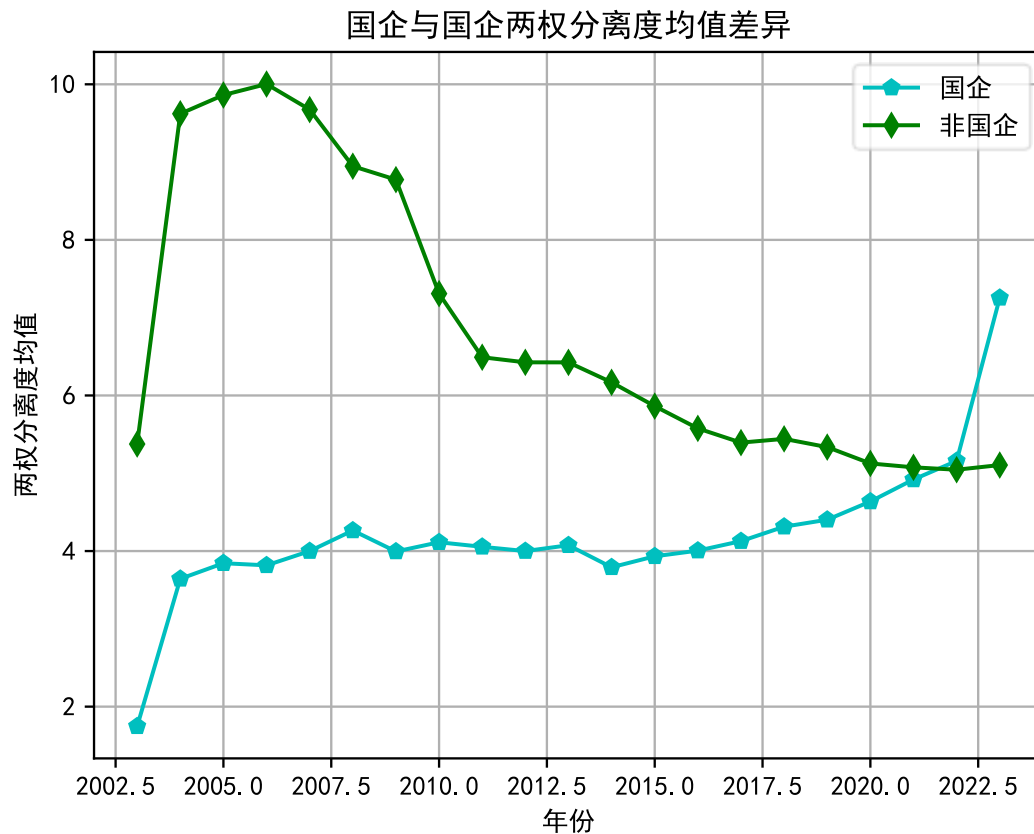
plt.legend()
plt.xlabel('年份')
plt.ylabel('两权分离度均值')
plt.title('国企与国企两权分离度均值差异')
plt.grid(True)

plt.show()

Text(0.5, 0, '年份')

Text(0, 0.5, '两权分离度均值')
```

```
Text(0.5, 1.0, '国企与国企两权分离度均值差异')
```



2.5.2.2 机构持股信息可视化

```
import pandas as pd
import matplotlib.pyplot as plt
from matplotlib import gridspec

#作图数据
status_counts = df['Status'].value_counts()
labels = status_counts.index.tolist()
sizes = status_counts.values.tolist()

annual_proportion_mean = df.groupby('year')['Proportion'].mean().reset_index()

annual_median = df.groupby('year')['Shares'].median().reset_index()

annual_sum=df.groupby('year')['HolderNumber'].sum().reset_index()

fig = plt.figure(figsize=(12,12))
grid = gridspec.GridSpec(2, 2, height_ratios=[2, 1], width_ratios=[1, 1])
```

```

# 绘制机构增减持饼图
# 0: 维持; 1: 增持; 2: 减持; 3: 新进
ax1 = fig.add_subplot(grid[0, 0])
ax1.pie(sizes, labels=labels, autopct='%1.1f%%', startangle=90)
ax1.axis('equal') # 确保饼图是圆的
ax1.set_title('机构增减持图')

# 绘制机构持股数年度比例均值折线图
ax2 = fig.add_subplot(grid[0, 1])
ax2.plot(annual_proportion_mean['year'], annual_proportion_mean['Proportion'], marker='o', linestyle='-', color='blue')
ax2.set_title('机构持股比例年度均值折线图')
ax2.set_xlabel('年份')
ax2.set_ylabel('机构持股比例年度均值')
ax2.grid(True)

# 绘制机构持股数年度中位数柱状图
ax3 = fig.add_subplot(grid[1, 0])
ax3.bar(annual_median['year'], annual_median['Shares'], color='c')
ax3.set_title('机构持股数年度中位数柱状图')
ax3.set_xlabel('年份')
ax3.set_ylabel('机构持股数年度中位数')
ax3.grid(True)

# 绘制机构持股数家数年度总和折线图
ax4 = fig.add_subplot(grid[1, 1])
ax4.plot(annual_sum['year'], annual_sum['HolderNumber'], marker='p', linestyle='-', color='g')
ax4.set_title('机构持股家数年度总和折线图')
ax4.set_xlabel('年份')
ax4.set_ylabel('机构持股数年度持股总和')
ax4.grid(True)

# 调整子图间距
plt.tight_layout()

plt.show()

([<matplotlib.patches.Wedge at 0x19a3ab3b200>,
 <matplotlib.patches.Wedge at 0x19a32cccc20>,
 <matplotlib.patches.Wedge at 0x19a38e19fa0>,
 <matplotlib.patches.Wedge at 0x19a33067d70>],
 [Text(-0.9640959812767127, 0.5296403863812623, '3.0'),
  Text(-0.1072154587189917, -1.094762460724552, '0.0'),
  Text(1.0928611753575848, -0.12511775012378865, '2.0'),
  Text(0.54534764312417, 0.9552988789582626, '1.0')],
 [Text(-0.5258705352418432, 0.28889475620796123, '34.0%'),

```

```
Text(-0.058481159301268196, -0.5971431603952102, '28.9%'),
Text(0.5961060956495916, -0.06824604552206652, '20.6%'),
Text(0.2974623507950018, 0.5210721157954159, '16.5%'))])

(-1.0999998785819671,
1.0999981517220176,
-1.0999996551377191,
1.0999999835779866)

Text(0.5, 1.0, '机构增减持图')

Text(0.5, 1.0, '机构持股比例年度均值折线图')

Text(0.5, 0, '年份')

Text(0, 0.5, '机构持股比例年度均值')

Text(0.5, 1.0, '机构持股数年度中位数柱状图')

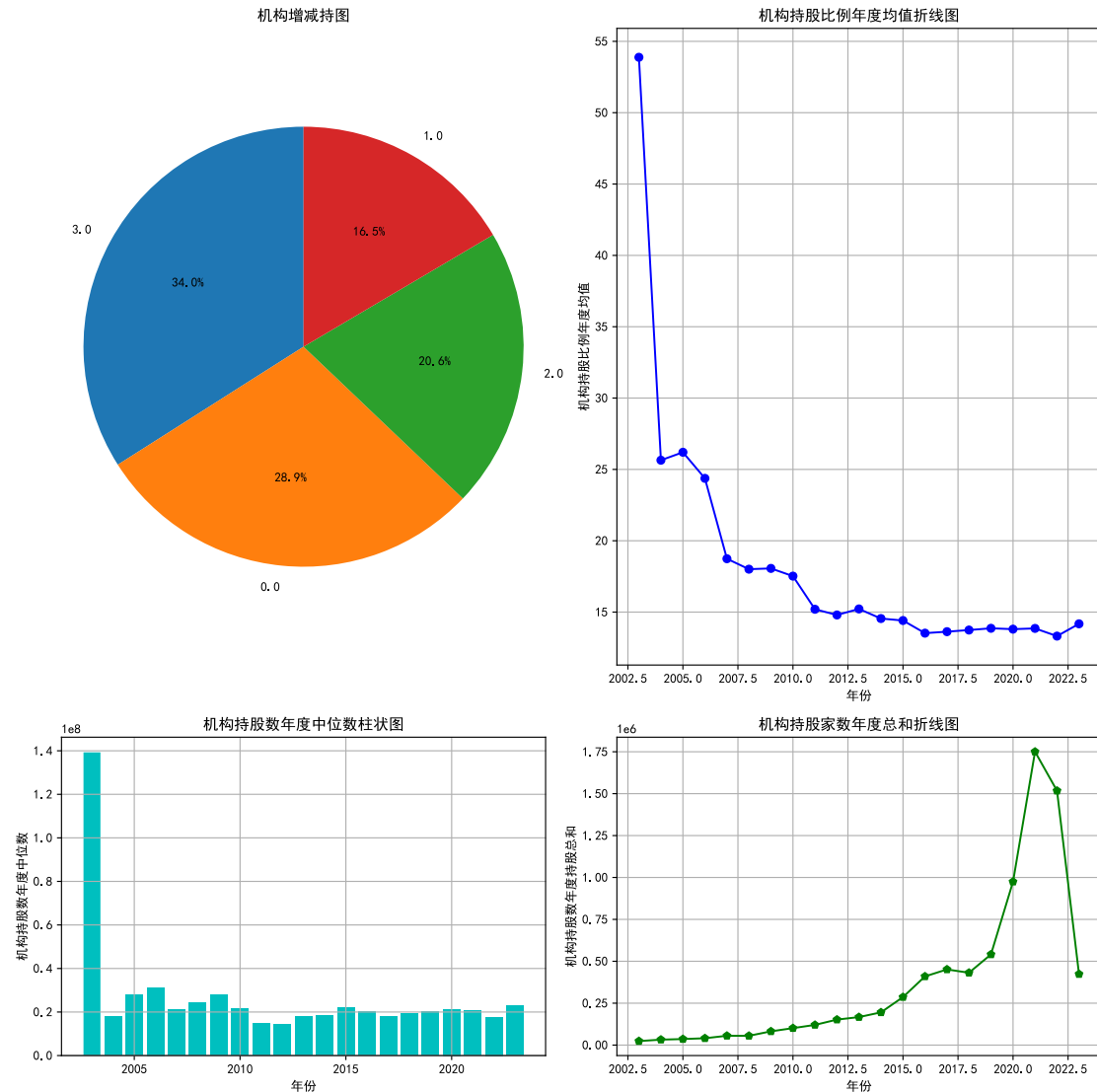
Text(0.5, 0, '年份')

Text(0, 0.5, '机构持股数年度中位数')

Text(0.5, 1.0, '机构持股家数年度总和折线图')

Text(0.5, 0, '年份')

Text(0, 0.5, '机构持股数年度持股总和')
```



2.5.2.3 分组指标探究

2.5.2.3.1 制造业与非制造业

#探寻制造业与非制造业第一大股东持股比例均值差异

```
import matplotlib.pyplot as plt
```

```
means = df.groupby(['year', 'manufacturingindustry'])['Shrcr1'].mean().reset_index()
```

```
yes_means = means[means['manufacturingindustry'] == 1]
```

```
no_means = means[means['manufacturingindustry'] == 0]
```

绘制制造业的数据

```
plt.plot(yes_means['year'], yes_means['Shrcr1'], label='Manufacturing', marker='o')
```

绘制非制造业的数据


```
plt.plot(no_means['year'], no_means['Shrcr1'], label='Non-Manufacturing', marker='s')

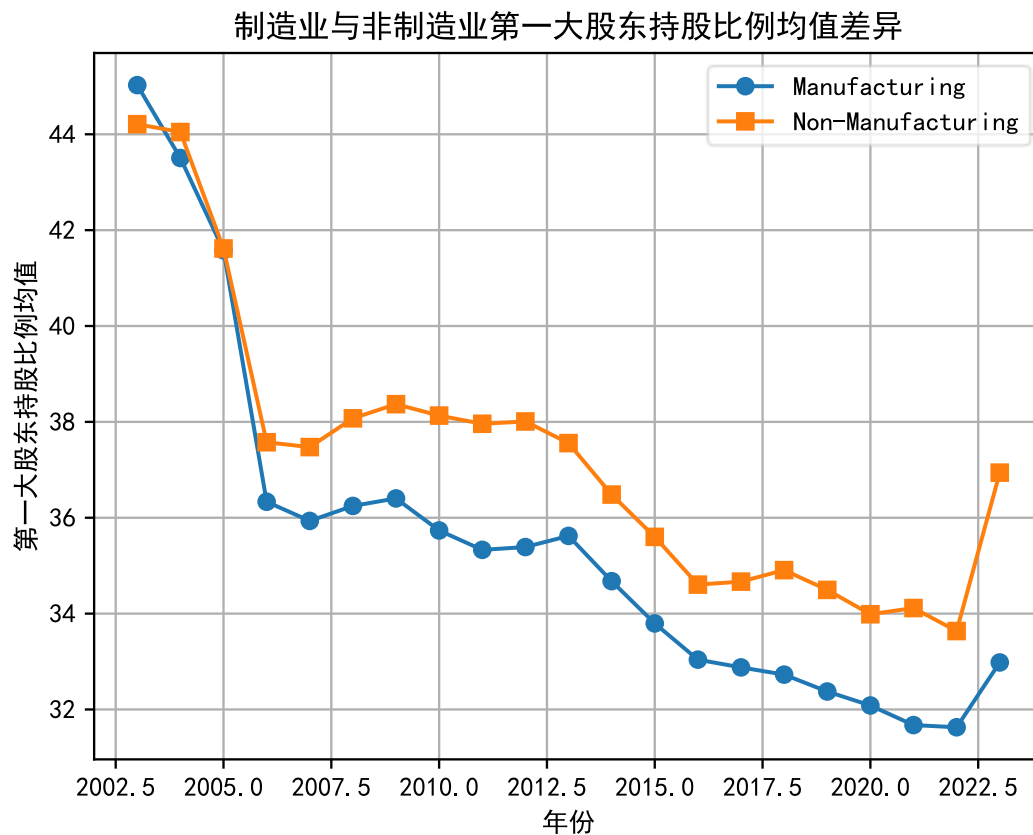
plt.legend()
plt.xlabel('年份')
plt.ylabel('第一大股东持股比例均值')
plt.title('制造业与非制造业第一大股东持股比例均值差异')
plt.grid(True)

plt.show()

Text(0.5, 0, '年份')

Text(0, 0.5, '第一大股东持股比例均值')

Text(0.5, 1.0, '制造业与非制造业第一大股东持股比例均值差异')
```



2.5.2.3.2 国企与非国企

#探寻国企与非国企第一大股东持股比例均值差异

```
import matplotlib.pyplot as plt
```

```
means = df.groupby(['year', 'PropertyRightsNature'])['Shrcr1'].mean().reset_index()
yes_means = means[means['PropertyRightsNature'] == 1]
```

```

no_means = means[means['PropertyRightsNature'] == 0]

plt.plot(yes_means['year'], yes_means['Shrcr1'], label='国企', marker='p')
plt.plot(no_means['year'], no_means['Shrcr1'], label='非国企', marker='s')

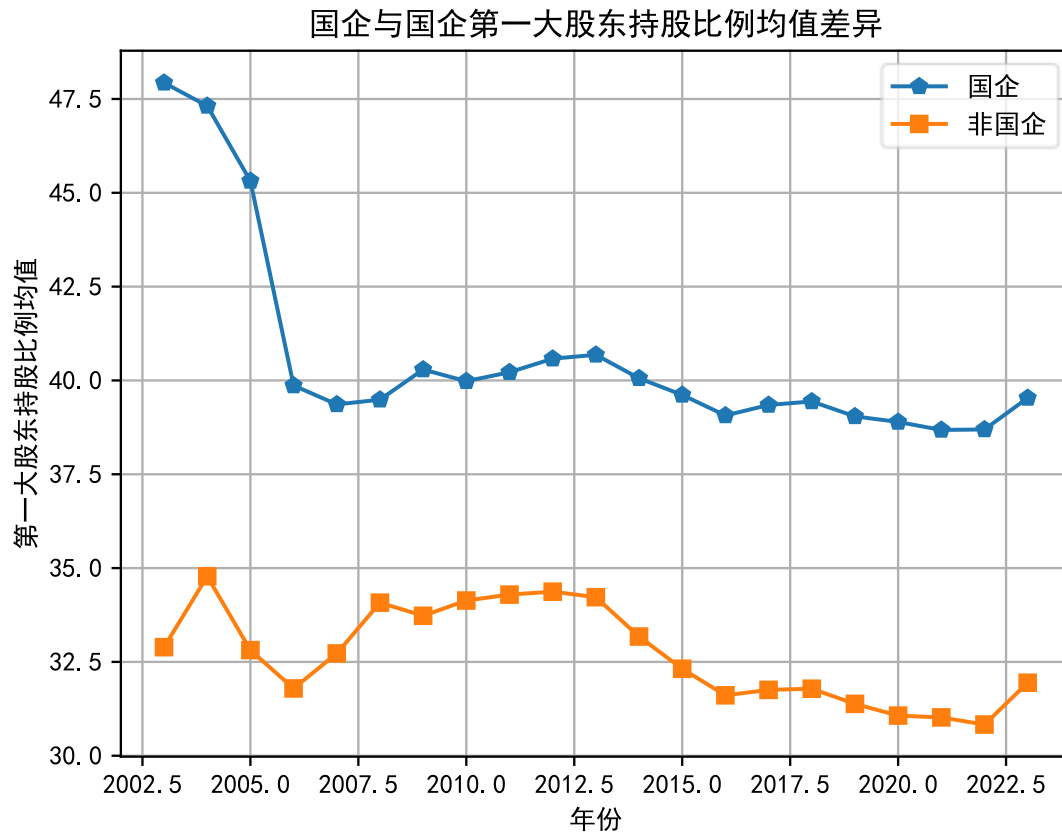
plt.legend()
plt.xlabel('年份')
plt.ylabel('第一大股东持股比例均值')
plt.title('国企与国企第一大股东持股比例均值差异')
plt.grid(True)
plt.show()

Text(0.5, 0, '年份')

Text(0, 0.5, '第一大股东持股比例均值')

Text(0.5, 1.0, '国企与国企第一大股东持股比例均值差异')

```



```

#探寻国企内持股比例均值差异
import matplotlib.pyplot as plt

```

```

#计算第一大股东比例均值
means1 = df.groupby(['year', 'PropertyRightsNature'])['Shrcr1'].mean().
reset_index()
yes_means1 = means1[means['PropertyRightsNature'] == 1]
#计算前三大股东比例均值
means2 = df.groupby(['year', 'PropertyRightsNature'])['Shrcr2'].mean().
reset_index()
yes_means2 = means2[means['PropertyRightsNature'] == 1]
#计算前五大股东比例均值
means3 = df.groupby(['year', 'PropertyRightsNature'])['Shrcr3'].mean().
reset_index()
yes_means3 = means3[means['PropertyRightsNature'] == 1]

plt.plot(yes_means1['year'], yes_means1['Shrcr1'], label='第一大股东持股
比例', marker='p')
plt.plot(yes_means2['year'], yes_means2['Shrcr2'], label='前三大股东持股
比例', marker='o')
plt.plot(yes_means3['year'], yes_means3['Shrcr3'], label='前五大股东持股
比例', marker='d')

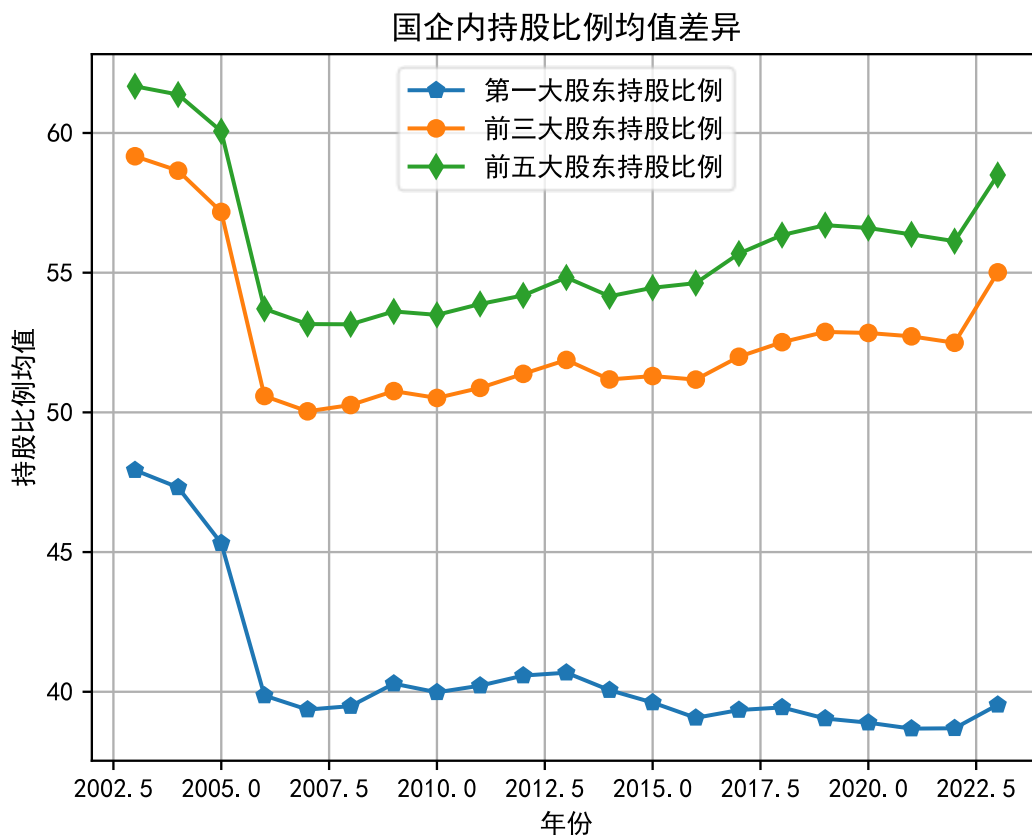
plt.legend()
plt.xlabel('年份')
plt.ylabel('持股比例均值')
plt.title('国企内持股比例均值差异')
plt.grid(True)
plt.show()

Text(0.5, 0, '年份')

Text(0, 0.5, '持股比例均值')

Text(0.5, 1.0, '国企内持股比例均值差异')

```



#探寻非国企内持股比例均值差异

```
import matplotlib.pyplot as plt
```

#计算第一大股东比例均值

```
means1 = df.groupby(['year', 'PropertyRightsNature'])['Shrcr1'].mean().reset_index()
```

```
no_means1 = means1[means['PropertyRightsNature'] == 0]
```

#计算前三大股东比例均值

```
means2 = df.groupby(['year', 'PropertyRightsNature'])['Shrcr2'].mean().reset_index()
```

```
no_means2 = means2[means['PropertyRightsNature'] == 0]
```

#计算前五大股东比例均值

```
means3 = df.groupby(['year', 'PropertyRightsNature'])['Shrcr3'].mean().reset_index()
```

```
no_means3 = means3[means['PropertyRightsNature'] == 0]
```

```
plt.plot(no_means1['year'], no_means1['Shrcr1'], label='非国企 1', marker='s')
```

```
plt.plot(no_means2['year'], no_means2['Shrcr2'], label='非国企 3', marker='o')
```

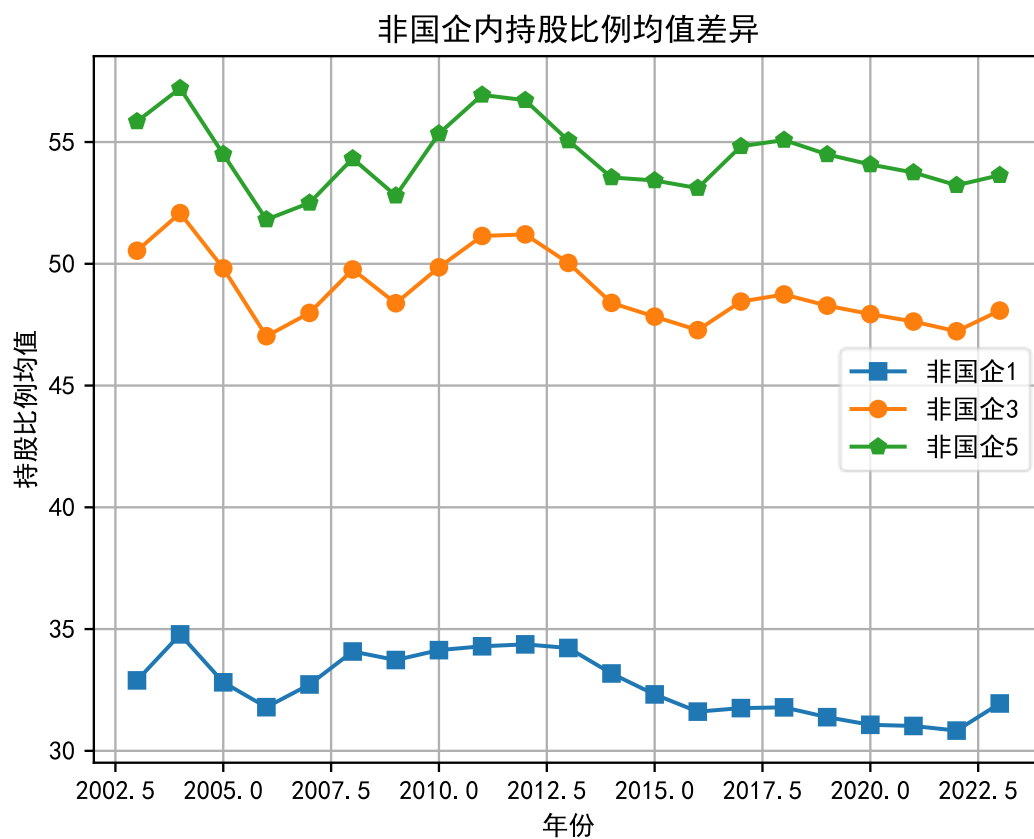
```
plt.plot(no_means3['year'], no_means3['Shrcr3'], label='非国企 5', marker='p')
```

```
plt.legend()
plt.xlabel('年份')
plt.ylabel('持股比例均值')
plt.title('非国企内持股比例均值差异')
plt.grid(True)
plt.show()

Text(0.5, 0, '年份')

Text(0, 0.5, '持股比例均值')

Text(0.5, 1.0, '非国企内持股比例均值差异')
```



3. 模型训练与评估

3.1 算法基本介绍

3.1.1 模型一、梯度提升树 GBDT 模型

梯度提升树（Gradient Boosting Decision Tree, GBDT）是一种集成学习方法，通过逐步构建并结合多个弱学习器（通常是决策树）来构建一个强大的预测模型。

核心思想是通过迭代的方式逐步减少残差（即预测值与真实值之间的差距）。在每一轮迭代中，模型会聚焦于当前模型预测结果与真实值之间残差的最大部分，然后训练一个新的决策树来拟合这部分残差。新树的输出被加到现有的模型预测上，从而逐步改善整体模型的预测能力。

在训练每棵决策树时，GBDT 采用了梯度下降的思想来最小化损失函数。具体来说，模型通过计算损失函数关于当前预测值的负梯度，确定残差的方向和大小，然后以该梯度作为指导构建决策树。这使得新树能够针对性地拟合残差，有效地减少整体损失。

3.1.2 模型二、Lasso 回归

Lasso 回归（Least Absolute Shrinkage and Selection Operator），也称为 L1 正则化线性回归，是一种数据挖掘方法，用于在常用的多元线性回归中添加惩罚函数，以压缩系数，从而达到精简模型的目的。

核心原理是在传统的线性回归损失函数中加入了一个 L1 正则化项（即参数的绝对值之和），通过正则化项对模型的参数进行惩罚，使得参数趋向于较小的值，同时也有助于产生稀疏模型，即某些特征的权重为零。

优势：

参数收缩与特征选择：通过 L1 正则化项，Lasso 回归可以将某些回归系数精确地压缩到 0，从而实现特征选择的目的。这使得模型更为简洁，减少了模型的复杂度。

防止过拟合：在拟合过于复杂的模型时，Lasso 回归通过正则化项对系数进行惩罚，有助于防止过拟合现象。

3.2 模型训练

```
df.columns
```

```
Index(['Stkcd', 'ShortName', 'Indcd', 'Indnme', 'F100801A', 'F100901A',  
      'F100902A', 'F100903A', 'F100904A', 'year', 'manufacturingindust  
ry',  
      'InstitutionID', 'Status', 'HolderNumber', 'Shares', 'Proportion',  
      'S0701a', 'S0702a', 'S0701b', 'S0703a', 'S0704a', 'S0705a', 'S07  
02b',
```

```

        'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNatur
e',
        'accper', 'A001212000', 'A001218000', 'A001000000', 'A002000000',
        'B001101000', 'B002000000', 'LargestHolderRate', 'TopTenHoldersR
ate',
        'EstablishDate', 'lnAssets', 'Lev', 'PPE_TA', 'Intass', 'OwnCon1
',
        'OwnCon2_10', 'SOE', 'ROA', 'Establish_year', 'Age',
        'Last_Year_Revenue', 'SalesGrowth'],
        dtype='object')

```

3.2.1 模型一

```

import pandas as pd
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.model_selection import train_test_split

X_value= ['Status', 'HolderNumber', 'Shares', 'Proportion', 'S0704a', '
S0705a', 'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNat
ure',
          'lnAssets', 'Lev', 'PPE_TA', 'Intass', 'OwnCon1', 'OwnCon2_10
', 'SOE', 'ROA', 'Establish_year', 'Age', 'SalesGrowth']
X1 = df[X_value]

y1 = df['F100901A']

X_train, X_test, y_train, y_test = train_test_split(
    X1, y1, test_size=0.2, random_state=42
)

gbreg = GradientBoostingRegressor(max_depth=2, n_estimators=3, learning
_rate=0.15, random_state=0)

gbreg.fit(X_train, y_train)

gbreg.predict(X_test[:5])

f"GradientBoostingRegressor 在训练集上的 R2:{gbreg.score(X_train, y_train):.3f}"
f"GradientBoostingRegressor 在测试集上的 R2:{gbreg.score(X_test, y_test):.
3f}"

GradientBoostingRegressor(learning_rate=0.15, max_depth=2, n_estimators
=3,
                           random_state=0)

```

```
array([2.0508443, 2.0508443, 2.0508443, 2.0508443, 2.0508443])
```

```
'GradientBoostingRegressor 在训练集上的 R2:0.327'
```

```
'GradientBoostingRegressor 在测试集上的 R2:0.221'
```

3.2.2 模型二

```
from sklearn.linear_model import Lasso
```

```
X_value= ['Status', 'HolderNumber', 'Shares', 'Proportion', 'S0704a', 'S0705a', 'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNature',
```

```
          'lnAssets', 'Lev', 'PPE_TA', 'Intass', 'OwnCon1', 'OwnCon2_10', 'SOE', 'ROA', 'Establish_year', 'Age', 'SalesGrowth']  
X2 = df[X_value]
```

```
y2 = df['F100901A']
```

```
X_train, X_test, y_train, y_test = train_test_split(X2, y2, test_size=0.2, random_state=42)
```

```
lasso = Lasso(alpha=0.95)
```

```
lasso .fit(X_train, y_train)
```

```
lasso .predict(X_test[:5])
```

```
lasso.score(X_train, y_train)
```

```
lasso.score(X_test, y_test)
```

```
lasso = Lasso(alpha=0.95).fit(X_train, y_train)
```

```
f"Training set score: {lasso.score(X_train, y_train):.3f}"
```

```
f"Test set score: {lasso.score(X_test, y_test):.3f}"
```

```
Lasso(alpha=0.95)
```

```
array([1.9151229 , 2.15419853, 1.81000124, 2.23009166, 2.10087588])
```

```
0.008929994856189283
```

```
0.011396955874669712
```

```
'Training set score: 0.009'
```

```
'Test set score: 0.011'
```


3.3 模型评估

3.3.1 模型一

```
from sklearn.metrics import (
    median_absolute_error,
    mean_absolute_error,
    mean_squared_error,
    r2_score,
)
from sklearn.model_selection import train_test_split
from sklearn.model_selection import RandomizedSearchCV
from IPython.core.interactiveshell import InteractiveShell

def scoring_reg(y_true, y_pred):
    print("-" * 10, "\n")
    print("median_absolute_error:", median_absolute_error(y_true, y_pred))
    print("mean_absolute_error:", mean_absolute_error(y_true, y_pred))
    print("mean_squared_error:", mean_squared_error(y_true, y_pred))
    print("r2_score:", r2_score(y_true, y_pred))
```

3.3.2 模型二

```
def myscoring(y_true, y_pred):
    print("-" * 10, "\n")
    print("median_absolute_error:", median_absolute_error(y_true, y_pred))
    print("mean_absolute_error:", mean_absolute_error(y_true, y_pred))
    print("mean_squared_error:", mean_squared_error(y_true, y_pred))
    print("r2_score:", r2_score(y_true, y_pred))
```

4. 模型调参

4.0.1 模型一

```
from scipy.stats import uniform
from scipy.stats import randint

params = {
    "max_depth": [3, 4, 5, 6, 7, 8],
    "min_samples_leaf": uniform(0.0001, 0.3),
    "learning_rate": uniform(0.0001, 0.5),
    "n_estimators": randint(10, 100)
}

gbdt_select = ['Status', 'HolderNumber', 'Shares', 'Proportion', 'S0704a', 'S0705a', 'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNature', 'lnAssets', 'Lev', 'PPE_TA', 'Intass', 'OwnCon1', 'Ow
```

```

nCon2_10', 'SOE', 'ROA', 'Establish_year', 'Age', 'SalesGrowth']
X1 = df[gbdt_select]

y1 = df['F100901A']

X_train, X_test, y_train, y_test = train_test_split(
    X1, y1, test_size=0.2, random_state=42
)

from sklearn.ensemble import GradientBoostingRegressor

rvestimator = GradientBoostingRegressor(random_state=0)
grid = grid = RandomizedSearchCV(
    estimator=rvestimator,
    param_distributions=params,
    n_iter=10,
    scoring="r2",
    n_jobs=-1,
    cv=5,
    random_state=0
)
grid.fit(X_train, y_train)

RandomizedSearchCV(cv=5, estimator=GradientBoostingRegressor(random_state=0),
                    n_jobs=-1,
                    param_distributions={'learning_rate': <scipy.stats._
distn_infrastructure.rv_continuous_frozen object at 0x0000019A345B24E0>,
                    'max_depth': [3, 4, 5, 6, 7, 8],
                    'min_samples_leaf': <scipy.stat
s._distn_infrastructure.rv_continuous_frozen object at 0x0000019A32FF44
A0>,
                    'n_estimators': <scipy.stats._d
istn_infrastructure.rv_discrete_frozen object at 0x0000019A34D98170>},
                    random_state=0, scoring='r2')

print("Best parameters: {}".format(grid.best_params_))
print("Best cross-validation score: {:.3f}".format(grid.best_score_))
print("Best estimator:\n{}".format(grid.best_estimator_))

Best parameters: {'learning_rate': 0.36041632736295837, 'max_depth': 4,
'min_samples_leaf': 0.04310598622271392, 'n_estimators': 42}
Best cross-validation score: 0.208
Best estimator:
GradientBoostingRegressor(learning_rate=0.36041632736295837, max_depth=
4,
                        min_samples_leaf=0.04310598622271392, n_estim
ators=42,
                        random_state=0)

```

4.0.2 模型二

```
param_lasso = {
    'alpha': [1e-4, 1e-3, 1e-2, 1e-1, 1, 10, 100], # 正则化强度
    'fit_intercept': [True, False] # 是否计算截距
}

lasso_select = ['Status', 'HolderNumber', 'Shares', 'Proportion', 'S070
4a', 'S0705a', 'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRig
htsNature',
                'lnAssets', 'Lev', 'PPE_TA', 'Intass', 'OwnCon1', 'Ow
nCon2_10', 'SOE', 'ROA', 'Establish_year', 'Age', 'SalesGrowth']
X2 = df[lasso_select]

y2 = df['F100901A']

X_train, X_test, y_train, y_test = train_test_split(
    X2, y2, test_size=0.2, random_state=42
)

from sklearn.linear_model import Lasso
lasso = Lasso()

la=la= RandomizedSearchCV(
    estimator=lasso,
    param_distributions=param_lasso,
    n_iter=10,
    scoring='neg_mean_squared_error',
    cv=5,
    random_state=42,
    n_jobs=-1
)

la.fit(X_train, y_train)

RandomizedSearchCV(cv=5, estimator=Lasso(), n_jobs=-1,
                   param_distributions={'alpha': [0.0001, 0.001, 0.01,
0.1, 1,
                                                10, 100],
                   'fit_intercept': [True, False]},
                   random_state=42, scoring='neg_mean_squared_error')

print("Best parameters: {}".format(la.best_params_))
print("Best cross-validation score: {:.3f}".format(la.best_score_))
print("Best estimator:\n{}".format(la.best_estimator_))

Best parameters: {'fit_intercept': True, 'alpha': 0.0001}
Best cross-validation score: -4.965
Best estimator:
Lasso(alpha=0.0001)
```

5. 模型再训练与评估

5.0.1 模型一

```
gbreg = GradientBoostingRegressor(  
    learning_rate=grid.best_params_["learning_rate"],  
    max_depth=grid.best_params_["max_depth"],  
    min_samples_leaf=grid.best_params_["min_samples_leaf"],  
    n_estimators=grid.best_params_["n_estimators"],  
    random_state=0  
)  
  
gbreg.fit(X_train, y_train)  
  
y_train_pred = gbreg.predict(X_train)  
y_test_pred = gbreg.predict(X_test)  
  
scoring_reg(y_train, y_train_pred)  
scoring_reg(y_test, y_test_pred)  
  
GradientBoostingRegressor(learning_rate=0.36041632736295837, max_depth=  
4,  
                           min_samples_leaf=0.04310598622271392, n_estim  
ators=42,  
                           random_state=0)  
-----  
  
median_absolute_error: 0.43055519778276263  
mean_absolute_error: 0.7104879538160958  
mean_squared_error: 4.208617129023834  
r2_score: 0.2097799142268003  
-----  
  
median_absolute_error: 0.4322560874143442  
mean_absolute_error: 0.7097417459687295  
mean_squared_error: 2.996913691821297  
r2_score: 0.2565593402969356
```

5.0.2 模型二

```
lareg = Lasso(  
    alpha=la.best_params_["alpha"],  
    fit_intercept=la.best_params_["fit_intercept"],  
    random_state=42  
)  
  
lareg.fit(X_train, y_train)  
  
y_train_pred = lareg.predict(X_train)  
y_test_pred = lareg.predict(X_test)
```

```

myscoring(y_train, y_train_pred)
myscoring(y_test, y_test_pred)

Lasso(alpha=0.0001, random_state=42)

-----

median_absolute_error: 0.6243878081397505
mean_absolute_error: 0.8640897125574449
mean_squared_error: 4.963253847374214
r2_score: 0.06808750695364008
-----

median_absolute_error: 0.6260719397293559
mean_absolute_error: 0.8583693732472447
mean_squared_error: 3.692201588441123
r2_score: 0.08408013478719556

```

6. 模型解释与应用

6.1 股权特征重要性

6.1.1 模型一

```

from sklearn.ensemble import GradientBoostingRegressor
import pandas as pd

X_value = ['Status', 'HolderNumber', 'Shares', 'Proportion', 'S0704a',
'S0705a', 'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNa
ture']
X1 = df[X_value]
y1 = df['F100901A']

gbdt_model = GradientBoostingRegressor()
gbdt_model.fit(X1, y1)

importance_gbdt = gbdt_model.feature_importances_

import matplotlib.pyplot as plt

indices = np.argsort(importance_gbdt)[::-1] # 索引从大到小排序
names = [X_value[i] for i in indices] # 获取对应的特征名称

plt.figure(figsize=(10, 5))
plt.bar(range(10), importance_gbdt[indices[:10]])

```

```
plt.title('GBDT 特征重要性(前十)')
plt.xlabel('特征')
plt.ylabel('重要性')
plt.xticks(range(10), [X_value[i] for i in indices[:10]], rotation='vertical')
plt.tight_layout()
plt.show()
```

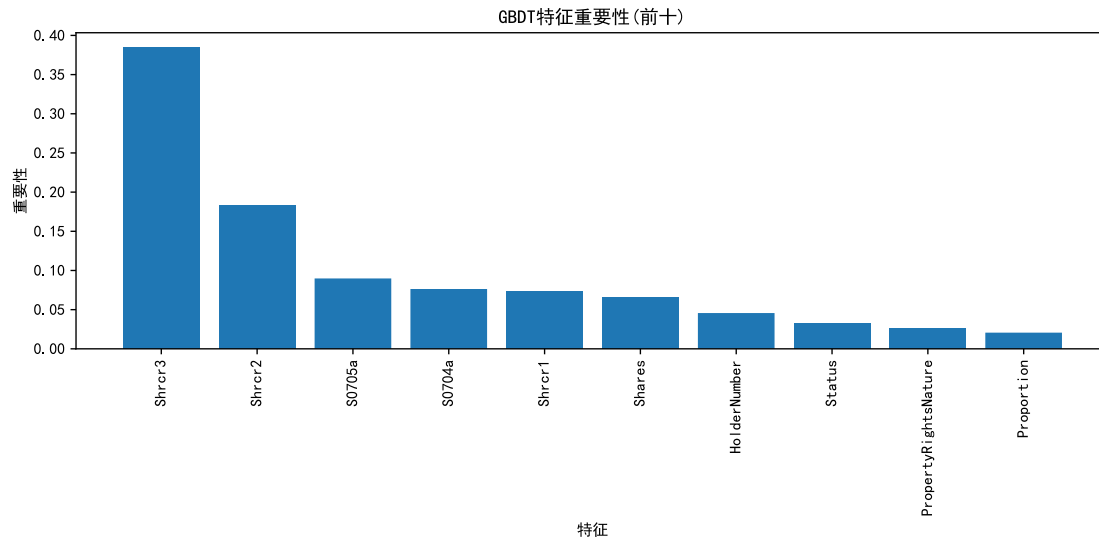
```
GradientBoostingRegressor()
```

```
Text(0.5, 1.0, 'GBDT 特征重要性(前十)')
```

```
Text(0.5, 0, '特征')
```

```
Text(0, 0.5, '重要性')
```

```
([<matplotlib.axis.XTick at 0x19a330142c0>,
 <matplotlib.axis.XTick at 0x19a34faf350>,
 <matplotlib.axis.XTick at 0x19a34f25df0>,
 <matplotlib.axis.XTick at 0x19a34f1bb60>,
 <matplotlib.axis.XTick at 0x19a34ef7d70>,
 <matplotlib.axis.XTick at 0x19a38447920>,
 <matplotlib.axis.XTick at 0x19a38446ff0>,
 <matplotlib.axis.XTick at 0x19a38447680>,
 <matplotlib.axis.XTick at 0x19a384466f0>,
 <matplotlib.axis.XTick at 0x19a38445d30>],
 [Text(0, 0, 'Shrcr3'),
 Text(1, 0, 'Shrcr2'),
 Text(2, 0, 'S0705a'),
 Text(3, 0, 'S0704a'),
 Text(4, 0, 'Shrcr1'),
 Text(5, 0, 'Shares'),
 Text(6, 0, 'HolderNumber'),
 Text(7, 0, 'Status'),
 Text(8, 0, 'PropertyRightsNature'),
 Text(9, 0, 'Proportion')])
```



6.1.2 模型二

```
from sklearn.linear_model import Lasso
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np

X_value = ['Status', 'HolderNumber', 'Shares', 'Proportion', 'S0704a',
'S0705a', 'Seperation', 'Shrcr1', 'Shrcr2', 'Shrcr3', 'PropertyRightsNa
ture']
X2 = df[X_value]
y2 = df['F100901A']

lasso_model = Lasso(alpha=0.95)
lasso_model.fit(X2, y2)

# 获取特征重要性（系数）
coef = lasso_model.coef_

print('前十特征重要性:')
for f in range(min(10, len(coef))):
    print(f'{X_value[f]}: {coef[f]}')

Lasso(alpha=0.95)

前十特征重要性:
Status: -0.0
HolderNumber: 0.0004987237506142459
Shares: -6.74619606217727e-12
Proportion: -0.0
```

```
S0704a: -0.00022074640644622007
S0705a: -1.5185300711834792e-11
Seperation: -0.0
Shrcr1: -0.004225190556922066
Shrcr2: -0.004841190641558022
Shrcr3: -0.0
```

7. 研究结论

7.1 股权特征典型事实

7.1.1 两权分离度

(1) 总体上，在 2003-2013 时间段呈现稳健上升特征，但在 2014-2022 时间段则是逐渐减少，2022 年后又再次上升。

(2) 国企与非国企之间，两者则是相反的差异，国企总体在提高两权分离度，而非国企总体在降低两权分离度。

7.1.2 机构投资者

(1) 在机构增减股份方面，选择维持与新进比例最多，而选择增持方面，占比最低。

(2) 每年度的机构持股数中位数分布上，除一开始的极端情况，总体分布较为均匀。

(3) 机构持股比例年度均值上呈现逐步下降的趋势，而机构持股数年度持股总和上先逐步上升，2018-2021 前后，持股总和快速上升，但与此同时，也以 2021 前后为转折点，持股总和又呈现出急剧下降的特点。

7.1.3 股权集中度

(1) 国企与非国企

在国企与非国企第一大股东持股比例均值上，国企均值明显大于非国企。趋势上，国企大体呈现下降趋势，最后在 2022 年下半年上升，非国企则是上升，下降，上升。

在第一，前三，前五持股比例之和均值的变化趋势:国企方面，三者基本呈现出先快速下降，但在除了第一大股东在 2006-2022 区间上是缓和下降的趋势，其余两者都是逐步上升。在 20022 年下半年，出现了较为快速的上升；非国企方面，前三，前五持股比例之和的均值波动比第一大股东更为明显，第一大股东均值变化较为缓和。

(2) 制造业与非制造业

在制造业与非制造业第一大股东持股比例均值上，2006 左右，两者从几乎没有差异发展到差异逐步拉开，但与此同时，两者总体都呈现出先下降，后在 2022 年左右又上升的趋势。

7.2 股权特征重要性

在 GBDT 中，Shrcr2 - 公司前 3 位大股东持股比例之和，与 Shrcr3 - 公司前 5 位大股东持股比例之和显著高于其他股权特征

报告承诺

本项目报告参考了：

[1]中国企业的专利行为:特征事实以及来自创新政策的影响（寇宗来 刘学悦）

[2]大数据应用对中国企业市场价值的影响——来自中国上市公司年报文本分析的证据（张叶青 陆瑶 李乐芸）

[3]数字技术创新对企业市场价值的影响研究（陶锋 朱盼 邱楚芝 王欣然）