

변수	통계/시각화	검정	특이사항
숫자 1	min,max,mean,std, median,사분위수  sns.histplot(data=df, x='v', kde=True) sns.boxplot(data=df, x='v') sns.kdeplot(data=df, x='v') sns.displot(data=df, x='v')	<정규분포(정규성) 검정> stats.shapiro(df.v) stats.anderson(df.v)  귀무가설: 정규분포이다.	<이상치 판단/처리> ① 사분위수 이용 ITQ=Q3-Q1 low=Q1-1.5*ITQ, //high=Q3+1.5ITQ df2=df[(df.v <high) & (df.v >low)] ② <u>zscore</u> z=np.abs(df.v -df.v.mean())/df.v.std() df2=df[z <2.5]
범주 1	범주별 빈도수(value_counts)/ 범주별 비율 sns.barplot(data=df, x='V') sns.countplot(data=df, x='V') * p=df.V.value_counts() plt.pie(p.values, label=p.index, autopc='%0.2f%%')		
범주(종속) 숫자(독립)	V_1=df[df.V ==0] V_2=df[df.V ==1] sns.kdeplot(x='v', data=V_1) sns.kdeplot(x='v', data=V_2)  sns.histplot(x='v', data=df, hue=V )	<로지스틱 회귀모형을 통한 검정> 회귀계수 0인가 import statsmodels.api as sm model=sm.Logit(df.V, df.v) result=model.fit() 귀무: 해당 변수의 회귀계수가 0이다 --독립과 종속이 관련이 없다.	
숫자 1,2,3...	sns.heatmap(df, vmin=-1, vmax=1, annot=True)	<상관관계> df.corr(numeric_only=True) 판다스의 .corr() 함수	
범주2	교차표 cr=pd.crosstab(df.V1, df.V2) cr=pd.crosstab(df.V1, df.V2, normalize='index') : normalize='index' /'columns'/'all'  cr.plot.bar(stacked=True) sns.countplot(data=df, x=v, hue='X')	<독립성 검정 x2(카이제곱)> cr=pd.crosstab(df.V1, df.V2) #교차표 stats.ch2_contingency(cr)  귀무: 두 변수가 독립적이다(연관성 없음)	카이제곱 검정: 예측값과 관측값의 차이 측정 ① 적합도(선호도) 검정: 변수1 ② 독립성 검정: 하나의 모집단에 범주형 2개 :범주변수들 간 어떤 관계가 있는지 ③ 동질성 검정: 2개 이상 모집단 2개 이상

범주(독립) 숫자(종속)	<p>&lt;정규성 검정&gt;  V_1=df[df.V ==0]['v']  V_2=df[df.V ==1]['v']  sns.histplot(V_1, kde=True)  sns.histplot(V_2, kde=True)</p>	<p>&lt;단일표본 T검정(one-sample T test)&gt;  하나의 집단평균이 모평균과 동일한지 검정(신뢰구간)  stats.ttest_1samp(V_1, df.1.mean())  stats.ttest_1samp(V_2, df.V.mean()) #귀무: 모평균과 동일</p> <p>&lt;정규성 검정&gt; --표준정규분포를 갖는지  stats.shapiro(V_1)  stats.shapiro(V_2) # 귀무: 등분산을 가진다</p> <p>&lt;등분산 검정&gt; --같은 분산을 갖는지 #귀무: 갖는다  stats.bartlett(V_1,V_2) # 표본크기 크고 정규성 가진 데이터  stats.levene(V_1,V_2) # 정규분포에 영향 안받음(중양값)  stats.fligner(V_1,V_2) # 정규분포 안따르는 데이터/중양값</p>	<p>신뢰구간 구하기  n=len(v)-1 #자유도  mn=np.mean(v) #표준편차  sm=stats.sem(v) #표준오차  stats.t.interval(0.95, n, mn, sm)  --&gt; 95% 신뢰구간 -1.96~1.96</p>
범주2(독립) 숫자1(종속)	<p>sns.barplot(x=V, y=v, data=df)  # V 범주들 모두 그래프에 표시됨</p>	<p>&lt;독립표본 T검정&gt; 모집단이 서로 다를때 두 평균의 차이  stats.ttest_ind(V_1,V_2, equal_var=False)  귀무: 유의미한 평균 차이가 없다.  t통계량 결과 -2&gt; 또는 &gt;2이면 '차이가 크다'.</p> <p>&lt;Mann-Whitney U test 검정&gt; 중양값 차이  표본수가 적을때 비모수적인 분석// 평균비교X, 순위합검정  stats.mannwhitneyu(V_1,V_2)  귀무: 두 집단의 중양값 차이가 없다.</p>	
범주3(독립) 숫자1(종속)	<p>V_1=df[df.V ==0]['v']  V_2=df[df.V ==1]['v']  V_3=df[df.V ==2]['v']  sns.barplot(x='V', y='v', data=df)</p>	<p>&lt;F검정(ANOVA)&gt; 신뢰구간 검정  여러 집단의 분산을 통한 집단의 차이가 있는지  stats.f_oneway('V_1','V_2','V_3')  귀무: 차이가 없다  f통계량 결과 2-3이상의 값을 가지면 차이가 있다</p>	