

# Topic1: Data Structure and Methods of Series

! pip install plotly | import numpy as np / scipy as sp / pandas as pd

**Data Structure** Series and DataFrame (DF contains multiple columns of Series)

- sales[0] / sales.loc["Jan"] ← Jan's value. - logic selection. m=sales.median(), mask=sales>m, sales[mask] 筛选出m的行

- NaN / np.nan. - df.SoldFlag / df[ "SoldFlag" ] (第ith Series) - 按Ctrl+Space会出现可选择的方法

- Methods : flag.max() - Attributes : flag.index / .real number

reduce to a value

① Operators ② Aggregate methods: .mean(), .max(), .sum(), .is\_monotonic() ③ Conversion methods: .to\_datetime() other format

④ Manipulation methods: .sort\_values(), drop\_duplicates return Series with same index ⑤ Indexing and accessor methods: .loc., .iloc. or scalars return Series

⑥ Plotting: - plot ⑦ Transformation method: .unstack(), reset\_index(), .agg(), .transform(), .pivot ⑧ Attribute: .index, .dtype.

Algebraic Operation: column-wisely +, -, \*, /, % 取余 modulus, @ matrix multiplication, \*\* power, <, >, &, |. -gt(c) 太长了

! list \* 2 就会成 [1,2] → [1,2,1,2], 而 Series \* 2 是 1.2 → 2.4 - 直接用 A+B, 只要一边没有值就会加到NaN, 但 A.add(B, fill\_value=0) 或 add(b)

Aggregation and Conversion: - city.quantile([.1, .8, .99]) 第一个最大的位置 - argmax(), count(), skew(), std(), var(), - agg("mean").is\_unique() ← 单值

- value\_count() → count 每个数的个数. 排分布. - A.astype("float32") 换数据类型 - history.index = pd.to\_datetime(history.index) 换成时间格式

Manipulation - A.groupby() / A.apply(def) - 一个效果. 但 apply 处理 - A.where(A.le150, other=0) 保留<50的. 否则为0 - Missing: .isna() / .isnull() / .dropna()

- fill() 填昨天的数据, bfill() 填明天的, interpolate() 前后平均 - clip(lower, upper) 所有数值直接用给定数字替换 - sort\_values(ascending=False) 大到小 / .sort\_index()

- drop\_duplicates(keep="subset") 本身留第一个 (last) 留最后一个 (first) 全部扔掉 - rank() 不动顺序下, 对值排序 0 40 3 20 2 1 (重复: 1) → 1, 2, 2, 2, 4; method='min' → 1, 2, 2, 4; 'dense' → 1, 2, 2, 3

- A.replace([1,2], "lowest"), (to\_replace=[4,1], value=[3,2]) / (to\_replace=[4:3, 1:2]) / (to\_replace=[2:, value='Suzanne', regex=True]): Suzy → Suzanne

- Binning: .cut(A, bins=4) / (A, [-1, 0, 1, 2, 3]) 返回每个所在区间: 0 [1,2] 1 [2,3]

return a series  
return data frame

Indexing Operation - S.rename(index={"0": "first"}), A.reset\_index(drop=True) 先去3-3行从0起, 再把原本的DFA. 但 .reset\_index()

- rename\_axis("columnn") - loc["A"] scalar / loc[[ "A" ]] series. - S.sort\_index().loc["G": "P"] 把从G到P的列取出来 (但若要pxx, 需写更远的终止)

- loc[min(1,1)].loc[gt3] / loc[gt3] - A.iloc[0] 值, .iloc[[0]] - 行, .iloc[1:-2] 第二行, 倒数第二行, [0:3]/[2:] - .head(10), .tail(10)

- A.sample(10) 抽10个出来, sample(frac=0.5, replace=True) - A.filter(like="K") 把 "K" 相关的行列出来, case sensitive. - A.reindex(["A", "B"]) A 3 B NaN

Date and Time - date.dt.day\_name() Sunday / Friday, dt.is\_month\_end() True / False, dt.strftime("%m/%d/%y"), loc[:="01/01/13"]

- Shifting temp.shift(1) 第一行变成NaN - Rolling data. rolling(3).mean() 表示本身 + 之前2个. - .cumsum() 累和和. cummax(),

- expanding() 从头开始加和, 和rolling()一样可选min\_periods 默认认为1. - resemble("Q"/"2M") 季度/每个月 (用于时间)

+ sum().agg(sum) 就是2M输出一个数, +.transform("sum") 就是把这2M的60天都转换成一个值.

Plotting - np.random.normal(10, 5, 500) 画出500个点 M=10, S=5 - ① matplotlib: plt.plot 是 line, .hist

- ② pandas: temp.plot.hist(), .box(), .kde(), .line(), - unstack() 会将最内层的变成多

用 pandas 的好处是可以进行前序的操作, 如 sector.value.value\_count().plot.barh()

s.agg(func=None, axis=0, *args, **kwargs)	Returns a scalar if func is a single aggregation function. Returns a series if a list of aggregations are passed to func.	r.apply(func, args=None, kwargs=None)	Apply custom aggregation function to rolling group.
s.all(axis=0, bool_only=None, skipna=True, level=None)	Returns True if every value is truthy. Otherwise False	r.corr(other, method='pearson')	Returns correlation coefficient for 'pearson', 'spearman', 'kendall', or a callable.
s.any(axis=0, bool_only=None, skipna=True, level=None)	Returns True if at least one value is truthy. Otherwise False	r.count(other, method='pearson')	Returns count of non NaN values.
s.autocorr(lag=1)	Returns Pearson correlation between s and shifted s	r.cov(other, min_periods=None)	Returns covariance.
s.corr(other, method='pearson')	Returns correlation coefficient for 'pearson', 'spearman', 'kendall', or a callable.	r.max(axis=None, skipna=None, level=None, numeric_only=None)	Returns maximum value.
s.cov(other, min_periods=None)	Returns covariance.	r.min(axis=None, skipna=None, level=None, numeric_only=None)	Returns minimum value.
s.max(axis=None, skipna=None, level=None, numeric_only=None)	Returns maximum value.	r.mean(axis=None, skipna=None, level=None, numeric_only=None)	Returns mean value.
s.min(axis=None, skipna=None, level=None, numeric_only=None)	Returns minimum value.	r.median(axis=None, skipna=None, level=None, numeric_only=None)	Returns median value.
s.mean(axis=None, skipna=None, level=None, numeric_only=None)	Returns mean value.	r.quantile(q=.5, interpolation='linear')	Returns 50% quantile by default. Note returns Series if q is a list.
s.median(axis=None, skipna=None, level=None, numeric_only=None)	Returns median value.	r.sem(axis=None, skipna=None, level=None, ddof=1, numeric_only=None)	Returns unbiased standard error of mean.
s.prod(axis=None, skipna=None, level=None, numeric_only=None, min_count=0)	Returns product of s values.	r.std(axis=None, skipna=None, level=None, ddof=1, numeric_only=None)	Returns sample standard deviation.
s.quantile(q=.5, interpolation='linear')	Returns 50% quantile by default. Note returns Series if q is a list.	r.var(axis=None, skipna=None, level=None, ddof=1, numeric_only=None)	Returns unbiased variance.
s.sem(axis=None, skipna=None, level=None, ddof=1, numeric_only=None)	Returns unbiased standard error of mean.	r.skew(axis=None, skipna=None, level=None, numeric_only=None)	Returns unbiased skew.
s.std(axis=None, skipna=None, level=None, ddof=1, numeric_only=None)	Returns sample standard deviation.	.date	Property with a series of Python datetime.date objects.
s.var(axis=None, skipna=None, level=None, ddof=1, numeric_only=None)	Returns unbiased variance.	.day	Property with a series of day of month.
s.kurtosis(axis=None, skipna=None, level=None, numeric_only=None)	Returns unbiased kurtosis.	.day_name(locale='en_us')	Return the string day of week.
s.nunique(dropna=True)	Returns count of unique items.	.dayofweek	Property with a series of date of week as number (0 is Monday).
s.count(level=None)	Returns count of non-missing items.	.dayofyear	Property with a series of day of the year.
s.size	Number of items in series. (Property)	.days_in_month	Property with a series of number of days in month.
s.is_unique	True if all values are unique	.daysinmonth	Property with a series of number of days in month.
s.is_monotonic	True if all values are increasing	.floor(freq=None, ambiguous=None, nonexistent=None)	Return floor according to offset alias in freq. The nonexistent parameter controls DST time issues.
s.is_monotonic_increasing	True if all values are increasing	.hour	Property with a series of hour of date.
s.is_monotonic_decreasing	True if all values are decreasing	.is_leap_year	Property with a series of booleans if date is leap year.
		.is_month_end	Property with a series of booleans if date is end of month.
		.is_month_start	Property with a series of booleans if date is start of month.
		.is_quarter_end	Property with a series of booleans if date is end of quarter.
		.is_quarter_start()	Property with a series of booleans if date is start of quarter.
		.is_year_end	Property with a series of booleans if date is end of year.

# Topic2: Working With DataFrames

**DataFrame** - df["xx"] 是 series 的格式 - .index(), .shape(), .columns(), .describe(), .axes(): 0 是 index, 1 是 column  
 $\text{df.sum}(\text{axis}=0) / (\text{axis}=\text{"index"})$ , df.sort\_index(by="x", ascending=False), df.set\_index("x"), df['A'].value\_count()

df["Season"] = [1 if t < 4 else 2 if t < 7 else 3 if t < 10 else 4 for t in df["Month"]]

df.rename(columns = lambda n: n + "o").head(), assign(New=lambda x: x.AEP \* 100), replace([1, 2], np.nan).head()  
 $\text{df1} + \text{df2}$  其中一个没有就是 NaN, .concat([df1, df2], axis=1) 0 就是上下拼接, 1 是左右, .append().append() .pct\_change()

**Aggregation:** df.agg(["count", "mean"], axis=0), df.apply(np.sum, axis=0)

**Missing & duplicated values:** .isna(), .isnull().sum() 每行. sum() 总的, df.loc[df.index[3]] = 3, .dropna() 去掉为行

**Filtering & subset selection** - df[["A", "B"]], loc 是左包含, iloc 右不包含的, df[df["AEP"] > 15000]

df.query("AEP > 15000 and COMED > 12000"), mask=df["AEP"] > 15000, df.query("@mask and COMED > 12000")

df.loc[df.AEP > 24000, df.columns[-3:]]

**Groupby Operation** - df.groupby("WeekDay")["AEP"].mean() 返回七个时间的 AEP 平均值.

df.groupby("WeekDay")["AEP"].transform(lambda x: (x - np.mean(x)) / np.std(x)) 对每个值, 做去相对 group 的 mean std.

df.groupby("A").filter(lambda x: x["B"].mean() > 3): 把 df 分成 M, N 组, 选其中 B 列的 mean > 3 的组 (M 的所有行)

**Pivoting:**

Year	New
1990	9
	10
1994	9
	20
	30
	40

sf.pivot_table(index="Year", columns="New", values="Price", aggfunc="mean")	new	0	1
	Year	1990	10
pd.crosstab(index=sf.Year, columns=sf.New, values=sf.Price, aggfunc="mean")	1994	30	40
	Year	1990	10
sf.groupby(["Year", "New"]).Price.mean().unstack()	1994	30	40
	Year	1990	10

**Data Assembly**

⇒ performance.melt(id\_vars=["Name", "ID"], value\_vars=["Python", "Statistics"], var\_name="Course", value\_name="Score")  
 $\Leftarrow pmelt.pivot_table(index=[\text{Name}, \text{ID}], columns="Course", values="Score")$

- transposing data: df.T

**Interactive Plotting** px.histogram(), scatter(), line()

- import plotly.graph\_objects as go. lay = go.Layout(), fig = go.Figure(). fig.show()

fig.add\_trace(), fig.update\_layout() fig=make\_subsets()

```

.is_year_start
.microsecond
.minute
.month
.month_name(locale='en_us')
.nanosecond
.normalize()
.quarter
.round(freq=None,
ambiguous=None,
nonexistent=None)
.second
.strftime(date_format)
.time
.timetz
.to_period(freq)
.to_pydatetime()
.tz
.tz_convert(tz)
.tz_localize(tz,
ambiguous=None,
nonexistent=None)

s.plot(ax=None, style=None, logx=False,
logy=False, xticks=None, yticks=None,
xlim=None, ylim=None, xlabel=None,
ylabel=None, rot=None, fontsize=None,
colormap=None, table=False, **kwargs)

```

```

s.plot.bar(position=.5, color=None)
s.plot.bart(x=None, y=None, color=None)
s.plot.hist(bins=10)
s.plot.box()
s.plot.kde(bw_method='scott', ind=None)

s.plot.line(color=None)
s.plot.pie()

```

## Operator methods instead of operators

Pandas also provides methods for the standard operators.

Method	Operator	Description
s.add(s2)	s + s2	Adds series
s.radd(s2)	s2 + s	Adds series
s.sub(s2)	s - s2	Subtracts series
s.rsub(s2)	s2 - s	Subtracts series
s.mul(s2) s.multiply(s2)	s * s2	Multiplies series
s.rmul(s2)	s2 * s	Multiplies series
s.div(s2) s.truediv(s2)	s / s2	Divides series
	s.mod(s2)	s % s2 Modulo of series division
	s.rmod(s2)	s2 % s Modulo of series division
	s.floordiv(s2)	s // s2 Floor divides series
	s.rfloordiv(s2)	s2 // s Floor divides series
	s.pow(s2)	s ** s2 Exponential power of series
	s.rpow(s2)	s2 ** s Exponential power of series
	s.eq(s2)	s2 == s Elementwise equals of series
	s.ne(s2)	s2 != s Elementwise not equals of series
	s.gt(s2)	s > 2 Elementwise greater than or equals of series
	s.ge(s2)	s >= 2 Elementwise greater than or equals of series
	s.lt(s2)	s < 2 Elementwise less than or equals of series
	s.le(s2)	s <= 2 Elementwise less than or equals of series
	np.invert(s)	~s Elementwise inversion of boolean series (no pandas method)
	np.logical_and(s, s2)	s & s2 Elementwise logical and of boolean series (no pandas method).
	np.logical_or(s, s2)	s   s2 Elementwise logical or of boolean series (no pandas method).
'all'		Returns True if every value is truthy.
'any'		Returns True if any value is truthy.
'autocorr'		Returns Pearson correlation of series with shifted self. Can override lag as keyword argument (default is 1).
'corr'		Returns Pearson correlation of series with other series. Need to specify other.
'count'		Returns count of non-missing values.
'cov'		Returns covariance of series with other series. Need to specify other.
'dtype'		Type of the series.
'dtypes'		Type of the series.
'empty'		True if no values in series.
'hasnans'		True if missing values in series.

# Topic3: Reshaping Statistics with Python

Population and Sample	<code>-data['A'].sample(10, replace=True)</code>	
$\mu: \text{mean}()$ , $\sigma^2: \text{var(ddof=0)}$ , $\sigma: \text{std(ddof=0)}$ , $\text{size: shape[0]}$		$25\% \quad 50\% \quad 75\% + 1.5 \text{ IQR}$
Variable and Distribution	Numerical, Categorical.	$\text{IR: A.quantile(0.75)} - \text{A.quantile(0.25)}$
Kernel Density Estimation:	$\hat{f}(x; h) = \frac{1}{nh} \sum_i K\left(\frac{x-x_i}{h}\right)$	kernel function: $\int K_h(x) dx = 1$ Gaussian, Rectangular, Triangular, uniform
$h: \text{bandwidth}$ small $h: \text{small bias, large var (noise)}$ , model.score_sample(): output the log function of the densities.		
from sklearn.neighbors import KernelDensity	KernelDensity(kernel='gaussian', bandwidth=0.1),	
- 对于厚尾(fat tail)的 stats.t.fit(data)[0] 输出 df = 3.28. (+越大越靠近正态)		
Confidence interval and Hypothesis Testing:	$\bar{x} \rightarrow n, \frac{\bar{x}-\mu}{S/\sqrt{n}} \rightarrow Z$	$[\bar{x} - 2\frac{\epsilon}{\sqrt{n}}, \bar{x} + 2\frac{\epsilon}{\sqrt{n}}]$
normality: histogram to visual, stats.probplot(data, dist="norm") 是 Q-Q plot, S-W test: shapiro(data): $H_0: \text{normal}$ .		$p < 0.05$ 拒绝正态。
Kolmogorov-Smirnov Test: kstest(data, stats.norm.cdf)	$H_0: \text{from the same distribution}$	
$x_{\bar{b}} = .\text{mean}()$ , $S = .\text{std}(ddof=1)$ $n = .\text{shape}[0]$ $Z_{\text{hat}} = (X_{\bar{b}} - \mu) / (S / (n^{0.5}))$ , p-value = 1 - norm.cdf( $Z_{\text{hat}}, 0, 1$ )		
Association: corr(), scatter_matrix(), corr() 包含 "pearson", "Spearman", "Kendall"		$\text{C}_i \xrightarrow{\text{SIMILAR}}$
Pearson correlation: can find linear pattern, Spearman's Rank-Order correlation: association based on rank [1..1]		
Hoeffding Dependence Coefficient: sign is not interpretation, identify non-monotonicity, Distance Coefficient: 0 is indep [0..1]		
Extreme Value Theory: Kurtosis $E\left(\frac{X-\mu}{\sigma}\right)^4$		
Fisher-Tippett-Gnedenko Theorem iid. $M_n = \max\{X_1, \dots, X_n\} \approx \text{GEV}(n, \theta, \xi)$ (General Extreme Value Dist.)		
$F(s) = \begin{cases} \exp(-\exp(-s)) & \xi=0 \\ \exp(-\exp(-s))^{\frac{1}{1+\xi}} & \xi \neq 0 \end{cases}$	$\begin{cases} \mu \text{ is location parameter. shift left/right} \\ \sigma \text{ is scale parameter. the spread of the distribution} \\ \xi \text{ is extreme value index (EVI), shape parameter, determine shape of distribution} \end{cases}$	$\xi > 0: \text{Weibull Dist. nega. int.}$
gev.pdf(x, c=1, loc=0, scale=1)		
3.7.1. Fréchet Dist. positive infinite $(x > \mu - \frac{\theta}{\xi})$ Heavy tailed; $\xi=0$ , Gumbel Dist. pos & nega infinite. Light tailed.		
How to sample? - Block Maxima (BM) 每一段时间的最大(不会有关联性) - Point Above Threshold (POT) (取所有大值)		
C, loc, scale = gev.fit(data, loc=0, scale=0.7) set initial guess		
Extreme Score: tail prob of GEV. 有多少概率会更极端(损失超-0.05) $1 - \text{gev.cdf}(\text{data}, c, \text{loc}, \text{scale})$		
Rare event: extreme score is small 明显预测下大雨的概率几乎为0, 可还是下了		
Extreme Value at Risk: 如果算出 95% EVaR = -0.06069, 则只有 5% 的概率损失会大于 0.06069.		

# Topic4: Predictive Model

$SST = \sum (Y_i - \bar{Y})^2$ , $SSR = \sum (\hat{Y}_i - \bar{Y})^2$ $SSE = \sum (Y_i - \hat{Y}_i)^2$	$R^2_{\text{adj}} = \left[ R^2 - \frac{k}{n-k-1} \right] \left[ \frac{n-1}{n-k-1} \right] = 1 - \frac{SSE/(n-k-1)}{SST/(n-1)}$
$MSE = \frac{e_1^2 + \dots + e_n^2}{n-k-1}$	- overfitting / underfitting. 看 train/test 的 $R^2_{\text{adj}}$ , RMSE, confusion matrix 的比例
- Nonlinear: $\odot X_1^2, X_2^2, X_1 X_2, \odot \frac{X_1}{X_2}, \frac{X_2}{X_1}$	
import dcor. dcor.distance_correlation	$F(\text{model}) = \frac{SSR/k}{SSE/[n-(k+1)]}$ $F > 2.5$ then reject $H_0: \beta_1 = \dots = \beta_k = 0$
model.bse: sd of the parameter. model["pvalue"] = 2 * norm.cdf(-np.abs(model.params / model.bse))	
$P(Y=1 X) = \frac{e^{\beta_0 + \beta_1 X}}{1+e^{\beta_0 + \beta_1 X}}$	$PPV = P(\text{Real 1} \text{Predict 1})$
$\log\left(\frac{P}{P_F}\right) = \beta_0 + \beta_1 X$	$NPV = P(\text{Real 0} \text{Predict 0})$

$$\begin{aligned} \text{Sensitivity} &= P(\text{Predict 1}|\text{Real 1}) = TPR, \text{ True Positive Rate, TPR} \\ \text{Specificity} &= P(\text{Predict 0}|\text{Real 0}) \\ 1 - \text{Specificity} &= P(\text{Predict 1}|\text{Real 0}) = FPR, \text{ False Positive Rate, FPR} \end{aligned}$$

	realY=0	realY=1	Measure
predictedY=0	True Negative TN	False Negative FN	Negative predictive value(NPV) $\frac{TN}{TN+FN}$
predictedY=1	False Positive FP	True Positive TP	Positive predictive value(NPV) $\frac{TP}{TP+FP}$
Measure	specificity $\frac{TN}{FP+TN}$	sensitivity $\frac{TP}{TP+FN}$	Accuracy $\frac{TP+TN}{TP+FP+FN+TN}$