



品职教育
PZACADEMY.COM

Quantitative Methods CFA一级知识框架图



扫码查看电子版勘误



讲师：何 旋

www.pzacademy.com

Framework

Time Value Calculation	Rates and Returns
	The Time Value of Money in Finance
Probability & Descriptive Statistics	Statistical Measures of Asset Returns
	Probability Trees and Conditional Expectations
	Portfolio Mathematics
	Simulation Methods
Inferential statistics	Estimation and Inference
	Hypothesis Testing
	Parametric and Non-Parametric Tests of Independence
Linear Regression	Simple Linear Regression
Big Data	Introduction to Big Data Techniques



Module 1



RATES AND RETURNS

Interest rate

含义

- ① Required rate of return
- ② Discount rate
- ③ Opportunity cost

构成

- ① $(1 + \text{nominal risk-free rate}) = (1 + \text{real risk-free rate})(1 + \text{inflation premium})$
- ② Required interest rate on a security = nominal risk-free rate + default risk premium + liquidity risk premium + maturity risk premium

Rates of Return ★★

Holding Period Return

A single specified period of time: $R = \frac{(P_1 - P_0) + I_1}{P_0}$

Three-year holding period return: $R = [(1 + R_1) \times (1 + R_2) \times (1 + R_3)] - 1$

Mean Return

arithmetic mean: $\overline{R}_i = \frac{R_{i1} + R_{i2} + \dots + R_{iT-1} + R_{iT}}{T}$

geometric mean: $\overline{R}_{Gi} = \sqrt[T]{(1 + R_{i1}) \times (1 + R_{i2}) \times \dots \times (1 + R_{iT-1}) \times (1 + R_{iT})} - 1$

harmonic mean: $\overline{X}_H = \frac{n}{\sum_{i=1}^n (1/X_i)}$ with $X_i > 0$ for $i = 1, 2, \dots, n$

Mean Return	<p>结论:</p> <ul style="list-style-type: none"> ① harmonic mean \leq geometric mean \leq arithmetic mean; $H \times A = G^2$ ② the more disperse the observations, the greater the difference between the arithmetic and geometric means. ③ Arithmetic mean: one-period horizon; the sum of the deviations around the mean equals 0; sensitive to extreme values, or outliers. ④ Geometric mean: multi-period horizon; backward data; represents the growth rate or compound rate of return on an investment; excellent measure of past performance ⑤ Harmonic mean: useful in the presence of outliers; used most often when the data consist of rates and ratios, such as P/Es.
Trimmed Mean & Winsorized Mean	<p>结论:</p> <ul style="list-style-type: none"> ① Both seek to minimize the impact of outliers in a dataset. ② Trimmed mean 是移除极值后的均值, 即截尾均值. ③ Winsorized mean 是将极值用最近的观察值替代后的均值, 即缩尾均值.

<p>Money Weighted Return</p>	<p>定义: accounts for the money invested and provides the investor with information on the <i>actual return she earns</i> on her investment.</p> <p>计算:</p> <p>① 找到每一期的现金流 (每期的<i>时间间隔相同</i>) ;</p> <p>② MWRR = 计算<i>IRR</i></p> <p>性质:</p> <p>优点: 衡量 the <i>actual return the investor earns</i> on her investment.</p> <p>缺点: 1) 会受到现金流改变的影响, 所以不能衡量基金经理的业绩; 2) it <i>does not allow for a return comparison</i> between different individuals or different investment opportunities</p>
<p>Time-Weighted Returns</p>	<p>定义: measures the compound rate of growth of \$1 initially invested in the portfolio</p> <p>计算:</p> <p>① <i>Price the portfolio immediately prior to any significant addition or withdrawal</i> of funds;</p> <p>② <i>Calculate the holding period return on the portfolio for each subperiod</i>;</p> <p>③ <i>Link or compound holding period: 几何平均</i></p> <p>实务:</p> <p>① <i>The more frequent the valuation</i>, the <i>more accurate</i> the approximation. Daily valuation is commonplace.</p> <p>② Annualized time-weighted return : $R_{TW} = (1 + R_1) \times (1 + R_2) \times \dots \times (1 + R_{365}) - 1$</p> <p>性质 (优点) :</p> <ul style="list-style-type: none"> <i>not sensitive to the additions and withdrawals of funds</i>; the TWR is the <i>preferred performance measure for the evaluation of portfolios manager</i>

Annualized Return	① $R_{\text{annual}} = (1 + R_{\text{period}})^c - 1 \longrightarrow R_{\text{weekly}} = (1 + R_{\text{daily}})^5 - 1; R_{\text{weekly}} = (1 + R_{\text{annual}})^{1/52} - 1$ ② continuously compounded return from t to $t + 1$: $r_{t,t+1} = \ln(P_{t+1}/P_t) = \ln(1 + R_{t,t+1})$ ③ $r_{0,T} = r_{T-1,T} + r_{T-2,T-1} + \dots + r_{0,1}$	
Gross and Net Return	<ul style="list-style-type: none"> Gross return 扣除 <i>Trading expense; an appropriate measure for evaluating and comparing the investment skill of asset managers.</i> Net return 扣除 <i>all managerial and administrative expenses</i> (management expenses, custodial fees, or any other administrative expenses) 	
Pre-Tax and After-Tax Nominal Return	$R_{\text{After-Tax}} = R_{\text{Pre-Tax}} (1 - T)$	Gross return Net return Leveraged return After-tax nominal return After-tax real return <div>↓ 计算顺序</div>
Real Return	$(1 + \text{real return}) = (1 + \text{real risk-free rate})(1 + \text{risk premium})$	
Leveraged Return	$R_L = R_P + \frac{V_B}{V_E} (R_P - r_D)$	



Module 2



THE TIME VALUE OF MONEY IN FINANCE

计算PV ★★

Fixed-Income Instruments	Discount	$PV(\text{Discount Bond}) = FV_t / (1 + r)^t$
	Periodic Interest	$PV(\text{Coupon Bond}) = PMT_1 / (1 + r)^1 + PMT_2 / (1 + r)^2 + \dots + (PMT_N + FV_N) / (1 + r)^N$ $PV(\text{Perpetual Bond}) = PMT / r$
	Level Payments	$PV(\text{Annuity Instruments}) = A / (1 + r)^1 + A / (1 + r)^2 + \dots + A / (1 + r)^N$ $A = \frac{r(PV)}{1 - (1 + r)^{-t}}$
Equity Instruments	Constant Dividends	$PV_t = \frac{D_t}{r}$
	Constant Dividend Growth Rate	$PV_t = \frac{D_t(1 + g)}{r - g} = \frac{D_{t+1}}{r - g}$
	Changing Dividend Growth Rate (两阶段模型)	$PV = \sum_{i=1}^n \frac{D_t(1 + g_s)^i}{(1 + r)^i} + \frac{E(S_{t+n})}{(1 + r)^n}$ $E(S_{t+n}) = \frac{D_{t+n+1}}{r - g_l}$

计算I/Y和Growth ★★

Fixed-Income Instruments	<ul style="list-style-type: none"> Implied return(I/Y): discount rate or yield-to-maturity(YTM) YTM的假设: an investor expects to receive all promised cash flows; reinvest any cash received at the same YTM 计算: 按计算器, 已知PV、FV、PMT、N、求I/Y
Equity Instruments	<ul style="list-style-type: none"> 计算: $r = \frac{D_{t+1}}{PV_t} + g$; $g = r - \frac{D_{t+1}}{PV_t}$ Price-to-earnings Ratio & Forward Price-to-earnings Ratio: $\frac{PV_t}{E_t} = \frac{\frac{D_t}{E_t}(1+g)}{r-g}$ & $\frac{PV_t}{E_{t+1}} = \frac{\frac{D_{t+1}}{E_{t+1}}}{r-g}$ 根据P/E ratio和dividend payout ratio, 计算出Implied Growth, 与分析师预测的Growth进行比较: Implied Growth大于预测的Growth, 则判断高估; 反之低估。

Cash Flow Additivity


难点

Implied Forward Rates	$F_{1,1} = (1 + r_2)^2 / (1 + r_1) - 1$
Forward Exchange Rates	汇率形式是A/B, 则 $F = S_0 e^{(r_A - r_B) \times t}$

Option Pricing
(Construct a risk-free portfolio)

Call Option: buy h units of the underlying asset (S_0) and sell one call option (执行价 X)

$$C_1^u = S_1^u - X, \quad C_1^d = 0$$



$$V_1^u = h \times S_1^u - C_1^u, \text{ and } V_1^d = h \times S_1^d - 0.$$


$$V_1^u = V_1^d \rightarrow \text{求出 } h (\text{hedge ratio})$$

$$V_0 = V_1^u / (1+r) = V_1^d / (1+r)$$

$$V_0 = h \times S_0 - c_0 \rightarrow \text{求出 } c_0$$

Put Option: buy h units of the underlying asset and buy one put option (执行价 X)

$$P_1^u = 0, \quad P_1^d = X - S_1^d$$



$$V_1^u = h \times S_1^u + 0, \text{ and } V_1^d = h \times S_1^d + P_1^d.$$

$$V_1^u = V_1^d \rightarrow \text{求出 } h (\text{hedge ratio})$$

$$V_0 = V_1^u / (1+r) = V_1^d / (1+r)$$

$$V_0 = h \times S_0 + P_0 \rightarrow \text{求出 } P_0$$

Module 3

STATISTICAL MEASURES OF ASSET RETURNS

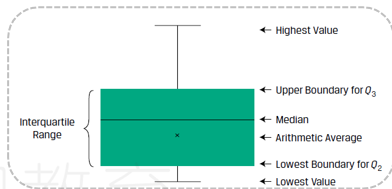
描述一组数据的基本特征

中心位置 ★★

Mean	<p>计算: arithmetic mean</p> <p>性质: the sum of the deviations around the mean equals 0</p> <p>缺点: <i>sensitive to extreme values, or outliers.</i></p> <p>对异常值的处理: ①Do nothing; ②Delete all the outliers(trimmed mean); ③Replace the outliers with another value(winsorized mean)</p>
Median	<p>计算: Odd numbered sample: the <i>$(n + 1)/2$ position</i></p> <p>Even-numbered sample: the mean of the values occupying <i>the $n/2$ and $(n + 2)/2$ positions</i></p> <p>优点: extreme values do not affect it; be useful <i>in describing data that follow a distribution that is not symmetric</i></p> <p>缺点: not use all the information about the size of the observations; focus only <i>on the relative position</i>; Calculating the median may also be <i>more complex.</i></p>
Mode	<p>计算: the <i>most frequently occurring value</i> in a distribution</p> <p><i>Unimodal, Bimodal and Trimodal</i></p>

Quantiles

Quantiles ★★	定义: Quartile / Quintile / Deciles / Percentile
	计算: $L_y = (n+1)y/100 \rightarrow$ linear interpolation
	性质: 比如 The third quartile > median
	IQR: $Q_3 - Q_1 \rightarrow$ Box and Whisker Plot



离散程度 ★★

Range	计算: Range = maximum value – minimum value
	Advantage: ease of computation Disadvantage: It cannot tell us how the data are distributed
MAD	$MAD = \frac{\sum_{i=1}^N X_i - \bar{X} }{n}$

<p>Sample Variance and Sample Standard Deviation</p>	<p>For population: $\sigma^2 = \frac{\sum_{i=1}^N (X_i - \mu)^2}{N}$ For sample: $s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$</p> <ul style="list-style-type: none"> • $MAD < \sigma$ • The gap between the arithmetic mean and the geometric mean: $\bar{X}_G \approx \bar{X} - \frac{s^2}{2}$ <ul style="list-style-type: none"> • the larger the variance of the sample, the wider the difference between the geometric mean and the arithmetic mean
<p>Target Downside Deviation</p>	<p>Target Semivariance = $\frac{\sum_{\text{for all } X_i \leq B} (X_i - B)^2}{n-1}$</p> <p>$n$ = the total number of observations in the sample</p>
<p>Coefficient of Variation</p>	<p>$CV = \frac{s_x}{\bar{X}} \times 100\%$</p> <ul style="list-style-type: none"> • relative dispersion • direct comparisons of dispersion across different datasets. • a scale-free measure (it has no units of measurement)

Skewness

★★★掌握性质

计算	$S_K = \left[\frac{n}{(n-1)(n-2)} \right] \frac{\sum_{i=1}^n (X_i - \bar{X})^3}{s^3} \approx \left(\frac{1}{n} \right) \frac{\sum_{i=1}^n (X_i - \bar{X})^3}{s^3}$ <p>a symmetrical distribution: skewness = 0</p>	
性质	Positive skewed	<ul style="list-style-type: none">• Mode < median < mean• right long fat tail• frequent small losses and a few extreme gains• investors should be attracted by positive skewness
	Negative skewed	<ul style="list-style-type: none">• Mode > median > mean• left long fat tail• frequent small gains and a few extreme losses.

kurtosis

★★ 掌握性质

Leptokurtic	<ul style="list-style-type: none">• Sample kurtosis > 3; Excess kurtosis > 0• 相同 σ, 尖峰肥尾 (more frequent extremely large deviations from the mean than a normal distribution.)
Platykurtic	Sample kurtosis < 3; Excess kurtosis < 0
Normal distribution	Sample kurtosis = 3; Excess kurtosis = 0

Covariance & Correlation



	计算	性质
Covariance	$s_{XY} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n-1}$	<ul style="list-style-type: none"> How two variables move together The covariance of X with itself is equal to the variance of X Positive covariance: the random variables vary in the same direction. Negative covariance: the variables vary in the opposite direction.
Correlation	$r_{XY} = \frac{s_{XY}}{s_X s_Y}$	<ul style="list-style-type: none"> Correlation measures the strength of linear relationship between two random variables Standardization of covariance, $-1 \leq r_{XY} \leq +1$ If $r=0$, this doesn't indicate independence, it indicates an absence of any linear. A positive correlation close to +1 indicates a strong positive linear relationship; $r=1 \rightarrow$ perfect linear relationship. A negative correlation close to -1 indicates a strong negative linear relationship; $r=-1 \rightarrow$ perfect inverse linear relationship.

Limitation of Correlation Analysis

Scatter Plot: a graph that shows the relationship between the observations for two data series in two dimensions

Limitation of Correlation Analysis

- Two variables can have a **strong nonlinear relation** and still have a **very low correlation**.
- Correlation also may be an **unreliable** measure when **outliers** are present.
- Correlation does **not imply causation**.
- **spurious correlation**
 - (1) correlation between two variables that reflects **chance relationships** in a particular data set
 - (2) correlation induced by a **calculation** that mixes each of two variables **with a third**
 - (3) correlation between two variables arising not from a direct relation between them but **from their relation to a third variable**.



Module 4



PROBABILITY TREES AND CONDITIONAL EXPECTATIONS

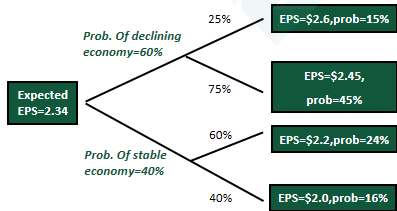
Expected Value and Variance ★

$$E(X) = P(X_1)X_1 + P(X_2)X_2 + \cdots + P(X_n)X_n = \sum_{i=1}^n P(X_i)x_i$$

$$\sigma^2(X) = E[X - E(X)]^2$$

$$\sigma^2(X) = P(X_1)[X_1 - E(X)]^2 + P(X_2)[X_2 - E(X)]^2 + \cdots + P(X_n)[X_n - E(X)]^2 = \sum_{i=1}^n P(X_i)[X_i - E(X)]^2$$

Probability Trees and Conditional Expectations ★★



会计计算Conditional Expectation:

$E(\text{EPS} | \text{declining economy})$

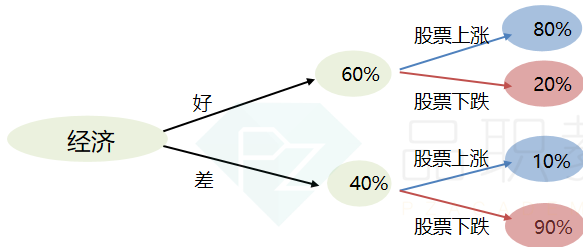
$E(\text{EPS} | \text{stable economy})$

$\sigma^2(\text{EPS} | \text{declining economy}) = P(2.6 | \text{declining economy}) \times [2.6 - E(\text{EPS} | \text{declining economy})]^2 + P(2.45 | \text{declining economy}) \times [2.45 - E(\text{EPS} | \text{declining economy})]^2$

$E(\text{EPS}) = E(\text{EPS} | \text{declining economy})P(\text{declining economy}) + E(\text{EPS} | \text{stable economy})P(\text{stable economy})$

Bayes' Formula

★★ 计算，用二叉树图形，不用记公式



注意：把非条件概率画在第一支

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$



Module 5



PORTFOLIO MATHEMATICS

Portfolio Expected Return and Variance

Covariance & Correlation	$\text{COV}(X,Y) = E[(X-E(X))(Y-E(Y))]$ $\rho_{XY} = \frac{\text{COV}(X,Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$
Independence ★★	<ul style="list-style-type: none"> • $P(AB)=P(A) \times P(B)$ • Independence $\rightarrow \rho=0$; 反之不对 • $E(XY) = E(X)E(Y)$ if X and Y are uncorrelated
Portfolio Expected Return	$E(r_p) = \sum_{i=1}^n w_i E(R_i)$
Portfolio Variance	<ul style="list-style-type: none"> • $\sigma_p^2 = w_1^2 \sigma_1^2 + w_2^2 \sigma_2^2 + 2\rho w_1 w_2 \sigma_1 \sigma_2$ • $\sigma_p^2 = w_1^2 \sigma_1^2 + w_2^2 \sigma_2^2 + w_3^2 \sigma_3^2 + 2\rho_{12} w_1 w_2 \sigma_1 \sigma_2 + 2\rho_{13} w_1 w_3 \sigma_1 \sigma_3 + 2\rho_{23} w_2 w_3 \sigma_2 \sigma_3$

Portfolio Risk Measures

1. Safety-First Ratio	$[E(R_p) - R_L] / \sigma_p \rightarrow \text{Sharpe Ratio} = [E(R_p) - R_f] / \sigma_p$ <p>Maximize SFR \Leftrightarrow Minimize $P(R_p < R_L)$: shortfall risk</p>
3. Stress testing/scenario analysis	estimating losses in extremely unfavorable combinations of events or scenarios .
4. Value at risk(VaR)	minimum value of losses expected over a specified time period at a given level of probability.



Module 6



SIMULATION METHODS

Lognormal Distribution and Continuous Compounding

类型	性质&计算
Lognormal distribution ★★	<ul style="list-style-type: none">• If $\ln X$ is normal, then X is lognormal.• Lognormal \rightarrow the price of asset; normal \rightarrow the return of asset• Right skewed; Bounded from below by zero (取值不能小于0)• Like the normal distribution, the lognormal distribution is completely described by two parameters (μ_L, σ_L^2).
Continuously Compounded Rates of Return ★★	<ul style="list-style-type: none">• Assume that the one-period continuously compounded returns (such as $r_{0,1}$) are i.i.d (independently and identically distributed).<ul style="list-style-type: none">• $r_{0,T}$ is approximately normal according to the central limit theorem $\rightarrow P_T$ is lognormal• $E(r_{0,T}) = E(r_{T-1,T}) + E(r_{T-2,T-1}) + \dots + E(r_{0,1}) = \mu T$• $\sigma^2(r_{0,T}) = \sigma^2 T$ (平方根原则)

Monte Carlo simulation ★

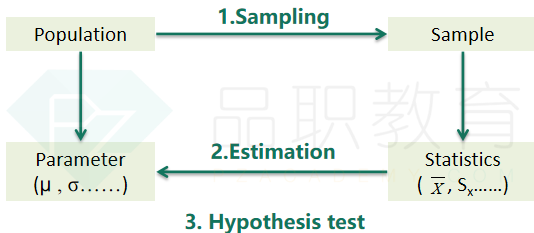
性质	
应用	<ul style="list-style-type: none">• is widely used to estimate risk and return in investment applications.• can be used to <i>price complex securities</i> for which no analytic expression is available.• examine the <i>model's sensitivity to a change in key assumptions</i>, such as mortgage-backed securities with complex embedded options.• address the sort of "what if " questions.
Limitations	<ul style="list-style-type: none">• <i>complex</i> and will assume a <i>parameter distribution</i>.• <i>complement to analytical methods</i> and <i>provides only statistical estimates</i>, not exact results. (Analytical methods provide more insight into cause-and-Effect relationships.)

Module 7,8,9

难点

ESTIMATION AND INFERENCE,
HYPOTHESIS TESTING, PARAMETRIC AND NON-PARAMETRIC
TESTS OF INDEPENDENCE

Framework



1. Sampling

Sampling methods ★	Probability sampling	<ul style="list-style-type: none">gives every member of the population an equal chance of being selected.The sample is representative of the population.
		<ol style="list-style-type: none">Simple random samplingSystematic sampling: select every kth memberStratified random sampling: draw simple random samples from each subpopulations in sizes proportional to the relative size of each subpopulations.Cluster Sampling: the population is divided into clusters, each of which is essentially a mini-representation of the entire populations. Then certain clusters are chosen as a whole.<ul style="list-style-type: none">Disadvantage: lower accuracy and less representativeAdvantage: time-efficient and cost-efficient
	non-probability sampling	<ul style="list-style-type: none">depends on other factors that is not probability considerationsMay generate a non-representative sample.
		<ol style="list-style-type: none">Convenience Sampling: an element is selected from the population based on whether or not it is accessible to a researcherJudgmental Sampling: select an item based on a researcher's knowledge and professional judgment.<ul style="list-style-type: none">Disadvantage: Sample selection could be affected by the bias of the researcher and might lead to skewed resultsAdvantage: allows researchers to go directly to the target population of interest.

Sample statistic 特点	<ul style="list-style-type: none"> sampling error of the mean= sample mean- population mean The sample statistic itself is a random variable 	
	Central Limit Theory ★★★	<ul style="list-style-type: none"> $n \geq 30 \rightarrow \text{sample mean} \sim N(\mu, \sigma^2/n)$ Standard error = σ / \sqrt{n} or s / \sqrt{n}



品职教育
P Z A C A D E M Y . C O M

Resampling ★

Bootstrap

- By treating the randomly drawn sample as if it were the population, we can simulate sampling from the population by sampling from the observed sample.
- the standard error of the sample mean:

$$S_{\bar{X}} = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}_b - \bar{\theta})^2}$$

Advantages:

- can be used to **find the standard error or construct confidence intervals** for the statistic of other population parameters, such as the median.
- is a **simple** but **powerful** method for any complicated estimators and particularly useful **when no analytical formula is available**.
- has potential advantages in **accuracy**.

Jackknife

- draw repeated samples while **leaving out** one observation at a time from the set, **without replacing it**
- For a **sample of size n** , jackknife usually requires **n repetitions**.

3. Hypothesis test

步骤: 检验 μ ★★

1. 提出假设

Two-tailed $H_0: \mu = \mu_0$ $H_a: \mu \neq \mu_0$

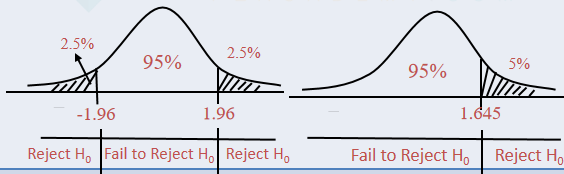
One-tailed $H_0: \mu \geq \mu_0$ $H_a: \mu < \mu_0$ or, $H_0: \mu \leq \mu_0$ $H_a: \mu > \mu_0$

H_0 is what we want to reject

2. 计算test statistic

$$\text{Test Statistic} = \frac{\bar{X} - \mu_0}{s / \sqrt{n}}$$

3. 画分布找到 critical value



4. 判断

Reject H_0 if $|\text{test statistic}| > \text{critical value} \rightarrow \text{*****}$ is significantly different from *****

Fail to reject H_0 if $|\text{test statistic}| < \text{critical value} \rightarrow \text{***}$ is not significantly different from ***

P – value ★★



P – value $< \alpha \rightarrow$ reject H_0

Type I error and Type II error ★★

Decision	True condition	
	H_0 ✓	H_0 ✗
Do not reject H_0	<u>Correct Decision</u>	<u>Incorrect Decision</u> Type II error
Reject H_0	<u>Incorrect Decision</u> Significance level = P (Type I error)	<u>Correct Decision</u> Power of test = 1 - P (Type II error)

1. Type I error $\uparrow \rightarrow$ Type II error \downarrow
2. Increase the Sample Size \rightarrow Type I error & Type II error \downarrow

Statistically and Economically Significant ★

- A strategy provides a statistically significant positive mean return. But the results may not be economically significant when we **account for transaction costs, taxes, and risk**.
- Even if we conclude that a strategy's results are economically meaningful, the economic logic of why the strategy might work **in the future** should be explored before implementing it.

Multiple Tests and Interpreting Significance ★

Multiple Tests	If you run 100 tests and use a 5% significance level , you get five false positives , on average.
False discovery rate (FDR)	<ul style="list-style-type: none">• false positive result: Type I error• The false discovery approach:<ul style="list-style-type: none">• adjusting the p-value when you have multiple tests: rank the p-values from the various tests, from lowest to highest, and then make the following comparison:$p(i) \leq \alpha \frac{\text{Rank of } i}{\text{Number of tests}}$• Repeat the comparison: k is determined by the highest ranked p(k) for which this is a true statement

其他检验 ★★★

Test type	Assumptions	H_0	Test-statistic	Critical value
Mean hypothesis testing	Normally distributed population, known population variance	$\mu=0$	$Z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$	$N(0,1)$
	Normally distributed population, unknown population variance	$\mu=0$	$t = \frac{\bar{x} - \mu_0}{s / \sqrt{n}}$	$t(n-1)$
	Independent populations, unknown population variances assumed equal	$\mu_1 - \mu_2 = 0$	$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}}$ <p><i>pooled estimator</i> of the common variance:</p> $s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$	$t(n_1 + n_2 - 2)$
	Independent populations, unknown population variances not assumed equal	$\mu_1 - \mu_2 = 0$	t	t
	Samples not independent, paired comparisons test	$\mu_d = 0$	$t = \frac{\bar{d}}{s_{\bar{d}}}$	$t(n-1)$
Variance hypothesis testing	Normally distributed population	$\sigma^2 = \sigma_0^2$	$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$	$\chi^2(n-1)$
	Two independent normally distributed populations	$\sigma_1^2 = \sigma_2^2$	$F = \frac{s_1^2}{s_2^2}$	$F(n_1 - 1, n_2 - 1)$

Parameter Tests and Non-parameter Tests ★

Parameter Tests	<ul style="list-style-type: none">• Rely on <i>assumptions</i> regarding the <i>distribution</i> of the population• Specific to population <i>parameters</i>.
Non-parameter Tests	<ul style="list-style-type: none">• examine quantities <i>other than population parameters</i> or where <i>assumptions of the parameters are not satisfied</i>. Nonparametric tests are used:<ul style="list-style-type: none">• The <i>assumptions</i> that support a parametric test <i>are not met</i>.• When there are <i>outliers</i>: influence the parametric statistics but not the nonparametric statistics.• When <i>data are ranks</i> (ordinal measurement scale) rather than values.• The <i>hypothesis does not involve the parameters</i> of the distribution, such as testing whether a variable is normally distributed.

Tests Concerning Correlation and Independence ★★

Tests Concerning Correlation			
Parametric Test of a Correlation	<ul style="list-style-type: none"> $H_0: \rho=0$ $t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}, df = n-2$ Two-tailed test Decision rule: reject H_0 if $+t_{critical} < t$, or $t < -t_{critical}$ 		
The Spearman Rank Correlation Coefficient	<ul style="list-style-type: none"> The population departs from normality → a test based on the Spearman rank correlation coefficient, r_s. Be calculated on the ranks of the two variables: <ul style="list-style-type: none"> Rank the observations. Calculate the difference, d_i, between the ranks for each pair of observations $r_s = 1 - \frac{6\sum_{i=1}^n d_i^2}{n(n^2-1)}$ 		
Test type	H_0	Parameter	Test-statistic
The Spearman Rank Correlation Coefficient	$r_s = 0$	$r_s = 1 - \frac{6\sum_{i=1}^n d_i^2}{n(n^2-1)}$	$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}, df = n-2$
Independence	independence	$E_{ij} = \frac{(\text{Total row } i) \times (\text{Total column } j)}{\text{Overall total}}$	$\chi^2 = \sum_{i=1}^m \frac{(O_{ij} - E_{ij})^2}{E_{ij}}, df=(r-1)(c-1)$

Module 10

SIMPLE LINEAR REGRESSION

建模



ANOVA Table分析



检验模型



预测

1.建模

$$Y_i = b_0 + b_1 X_i + \varepsilon_i, i = 1, \dots, n$$

- Y_i = **dependent variable**, explained variable, predicted variable
- X_i = **independent variable**, explanatory variable, predicting variable.

Assumption



ε_i 均值为0, 方差不变, 不相关的正态序列

- Linearity: The relationship between the ***Y and the X is linear.***
 - X must ***not be random (non-stochastic)***
 - The residuals are random. The residuals should not exhibit a pattern when plotted against the independent variable
- Homoskedasticity: The variance of the ***residuals is the same for all observations***
- Independence: ***The observations***, pairs of Y_s and X_s , ***are independent of one another.*** This implies the ***regression residuals are uncorrelated*** across observations
- Normality: The residuals be ***normally distributed***

Coefficient估计



解释:

- ***The intercept:*** the ***value of the dependent variable*** if the value of the ***independent variable is zero.***
- ***The slope:*** the change in the Y for ***a one-unit change in the X.***

OLS: the sum of the squared differences between the observations on Y_i and the corresponding estimated value \hat{Y}_i is minimized.

$$b_1 = \frac{Cov(X, Y)}{Var(X)}$$

$$b_0 = \bar{Y} - b_1 \bar{X}$$

计算

2. ANOVA Table分析 ★★★

	df	SS	MSS
Regression	k=1	RSS	MSR=RSS/k
Error	n-k-1	SSE	MSE=SSE/(n-k-1)
Total	n-1	SST	-



$$\left. \begin{aligned} SSE &= \sum_{i=1}^n \left(Y_i - \hat{Y}_i \right)^2 \\ SSR &= \sum_{i=1}^n \left(\hat{Y}_i - \bar{Y} \right)^2 \end{aligned} \right\} SST = \sum_{i=1}^n \left(Y_i - \bar{Y} \right)^2 = SSR + SSE$$

Coefficient Determination (R^2)	计算: $R^2 = \frac{RSS}{TSS} = 1 - \frac{SSE}{TSS}$
	$R^2 = r_{\hat{Y}\hat{Y}}^2$ (多元都成立) $R^2 = r_{XY}^2$ (一元)
	解释: R^2 of 0.90 indicates that the variation of the independent variable explains 90% of the variation in the dependent variable.
SEE	计算: $SEE = \sqrt{\frac{SSE}{n-k-1}} = \sqrt{MSE}$
	性质: <ul style="list-style-type: none"> The smaller the standard error, the better the fit. The SEE is the standard deviation of the error terms in the regression.

F-test

- $H_0: b_1 = b_2 = b_3 = \dots = b_k = 0$; H_a : at least one $b_j \neq 0$ ($j = 1$ to k)
- $F = \frac{MSR}{MSE} = \frac{RSS/k}{SSE/(n-k-1)} \sim F(k, n-k-1)$
- reject H_0 : if F (test-statistic) $> F_c$ (critical value)
- The F-statistic in regression analysis is **one sided**, with the **rejection region** on the **right side**.

总结: Measures of the goodness of the fit

- The **coefficient of determination** and the **F-statistic** are **relative** measures of fit
- The **standard error of the estimate** is an **absolute** measure

3. 检验模型：回归分析相当于抽样估计

考试时给定条件



	Coefficient	Standard deviation	t-statistic	p-value
Intercept	\hat{b}_0	$s_{\hat{b}_0}$?	0.18
Slope	\hat{b}_1		?	<0.001

Hypothesis Tests

★★★

- $H_0: b_1=0$ (没有特殊说明, 题目中假设检验都是检验是否为0)
- $t = \frac{\hat{b}_1 - B_1}{s_{\hat{b}_1}}$
- Decision rule: reject H_0 if $+t_{\text{critical}} < t$, or $t < -t_{\text{critical}}$
- Rejection of the null means that the slope coefficient is different from zero

Hypothesis Tests

Features of simple linear regression

★★★

- the **t-statistic** used to test whether the **slope coefficient is equal to zero** and the **t-statistic** to test whether the **pairwise correlation is zero** are the same value.
- the **F-distributed test statistic** & t-statistic used to test whether the slope coefficient is equal to zero: **$t^2 = F$**

Hypothesis Tests of the Intercept

$$t_{\text{intercept}} = \frac{\hat{b}_0 - B_0}{s_{\hat{b}_0}}$$

Indicator Variable or Dummy Variable



定义	take on only the <i>values 0 or 1</i> as the <i>independent variable</i>
解释 b_0, b_1	$RET_i = b_0 + b_1 \text{ EARN}_i + \varepsilon_i$ Y = monthly returns, RET over a 30-month period X = indicator variable, EARN, that takes on a value of 0 if there is no earnings announcement that month and 1 if there is an earnings announcement <ul style="list-style-type: none"> • The intercept (0.5629): the mean of the returns for <i>non-earnings-announcement months</i>. • The slope coefficient (1.2098): the <i>difference</i> in means of returns between earnings-announcement and non-announcement months
Hypothesis Tests	<ul style="list-style-type: none"> • Test whether the mean monthly return is the same for both the non-earnings-announcement months and the earnings-announcement months: $H_0: \mu_{RET_{\text{earnings}}} = \mu_{RET_{\text{non-earnings}}}$ <ul style="list-style-type: none"> • Reject H_0: there is difference in the mean RET for the earnings-announcement and non-earnings-announcements months

4.预测(Predicted Value of Y)

<p>The predicted value for Y</p>	<p>给定X代入计算Y: $\hat{Y}_f = \hat{b}_0 + \hat{b}_1 X_f$ 注意: in decimal \rightarrow 代入百分号; in percent \rightarrow 百分号去掉</p>
<p>Prediction Interval ★★</p>	<p> $\hat{Y}_f \pm t_{\text{critical for } \alpha/2} S_f$ $s_f = s_e \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}} \approx s_e$ </p> <p>The standard error of the forecast depends on: (<i>the smaller S_f depend on:</i>)</p> <ul style="list-style-type: none"> • <i>The better the fit</i> of the regression model, the <i>smaller the standard error of the estimate</i> (s_e) • <i>The larger the sample size</i> (n) in the regression estimation. • <i>The closer</i> the X_f is to the mean of the independent variable

Functional Forms

类型	性质
the Log-Lin model: $\ln Y_i = b_0 + b_1 X_i$	The slope coefficient: the <i>relative change in the dependent variable</i> for an <i>absolute change in the independent</i> variable.
the Lin-Log model: $Y_i = b_0 + b_1 \ln X_i$	The slope coefficient: the <i>absolute change in the dependent</i> variable for a <i>relative change in the independent</i> variable.
the Log-Log model: $\ln Y_i = b_0 + b_1 \ln X_i$	This model is useful in calculating elasticities because the slope coefficient is the <i>relative change in Y for a relative change in X</i> .

P Z A C A D E M Y . C O M

Module 11

INTRODUCTION TO BIG DATA TECHNIQUES

概念理解

Fintech

Fintech = technological innovation + financial services and products

Areas of fintech development

Analysis of large datasets	<ul style="list-style-type: none">• <i>Traditional data</i>• <i>Alternative data from non-traditional data sources: social media and sensor networks</i>
Analytical tools	Artificial intelligence (AI)

Big Data

1. Sources of Big Data ★★★

Traditional data → *corporate data*
financial markets

Non-traditional data (alternative data) →

<i>Individuals</i>	<ul style="list-style-type: none">• Social media• News, reviews• Web searches, personal data	<i>Unstructured</i> Volume: <i>growing</i> dramatically
<i>Business Processes</i>	<ul style="list-style-type: none">• Transaction data• <i>corporate exhaust</i> (corporate supply chain information, banking records, and retail point-of-sale scanner data)	<i>structured</i> data leading or real-time indicators of business performance
<i>Sensors</i>	<ul style="list-style-type: none">• <i>Satellites</i> and Geolocation• smart phones, cameras, RFID chips• <i>Internet of Things</i>	<i>Unstructured</i> Volume: <i>greater</i>

2. Characteristics of Big Data ★★★

<i>Volume</i>	The <i>amount</i> of data is very <i>large</i>
<i>Velocity</i>	<i>real-time</i> communication
<i>Variety</i>	<i>different sources</i> and in a <i>variety</i> of <i>formats</i> (<i>Structured, Semi-structured and Unstructured</i>)
<i>Veracity</i>	Determining the <i>credibility</i> and <i>reliability</i> of different data sources is an important part.

Artificial Intelligence & Machine Learning

1. Artificial Intelligence

Artificial Intelligence: capable of performing tasks that have traditionally required human intelligence

Development of AI

- **“Expert system”:** simulate the knowledge base and analytical abilities of human experts
- **Neural networks**

2. Machine Learning ★★★

Definition	computer programs that are able to “learn” how to complete tasks , improving their performance over time with experience
Terms	<ul style="list-style-type: none">• “inputs” (a set of variables or datasets) & “outputs” (the target data)• Algorithm: “learns” from the data; “black box” approaches• Training dataset and validation dataset (evaluation dataset)• Overfitting & underfitting<ul style="list-style-type: none">• Overfitting: the ML model learns the input and target dataset too precisely; be “over-trained” on the data and treats noise in the data as true parameters; too complex model; prediction errors using a different dataset• Underfitting: treat true parameters as if they are noise; too simplistic model; fail to fully discover patterns
Challenges	<ul style="list-style-type: none">• Still require human judgement• The data must be clean and free of biases and spurious data• Require sufficiently large amounts of data

Types of Machine Learning

Supervised learning

based on **labeled** (or identified) training data

- predict whether local stock market performance will be up, down, or flat

Unsupervised learning

be **not** given **labeled** data; **describe the data and their structure**

- group companies into peer groups based on their characteristics

Deep learning (AI advances) **neural networks**, with many hidden layers, to perform multistage, non-linear data



品职教育
P Z A C A D E M Y . C O M

Tackling Big Data With Data Science

computer science (including machine learning) + **extracting information from Big Data**

Data Processing Methods

Capture: how the data are **collected** (**Low-latency systems**) and **transformed into a format** for the analytical process

Curation: ensuring **data quality and accuracy**

Storage: how the data will be **recorded**, **archived**, and **accessed** (**structured or unstructured**)

Search: how to **query data**

Transfer: move from **the underlying data source or storage location** to the underlying **analytical tool**

Data Visualization → { Traditional structured data: using **tables**, **charts**, and **trends**
Non-traditional unstructured data: 3D charts, **heat maps**, **tree diagrams**, and **network graphs**, **"tag cloud"** and **"mind map"**

Applications	特点
Text Analytics	from large, unstructured text- or voice-based datasets ;
	Include automated information retrieval; lexical analysis, or the analysis of word frequency; identify indicators of future performance.
Natural Language Processing	the intersection of computer science , artificial intelligence , and linguistics ; analyze and interpret human language .
	Include translation, speech recognition, text mining, sentiment analysis, and topic analysis

*Thank
You!*

