

Note 2 Multivariate analysis

- Logic
 - On Note 1,
 - modeling one kind of stock's price from the raw data, normal distribution, student t distribution, and general t distribution, and we can validate it, KS test etc.
 - some index from the data, E, Var, f(X) < - moment estimate, interval estimate, MLE, 蒙特卡洛模拟
 - 当然这些，尤其是蒙特卡洛模拟都是建立在样本量足够大的情况下，样本量很小的情况下怎么解决？bootstrap
 - 但仅仅研究一只股票是不够的，我们最终要构建的就是一个portfolio，并且股票之间很可能是related的。这样我们就不能用之前的方法认为他们是iid的
 - On Note 2,
 - 我们将考虑多个股票可能服从的分布，joint distribution, 期望, 方差 (normal, t), 同时还有多元的蒙特卡洛模拟
 - 对于多元函数来讲，我们还有另外的方法来模拟股票收益的走势，线性回归
 - CAPM
 - Fama-French three-factor model
- Basic concepts
 - joint distribution:
 - If X_i s are discrete random variables, the joint distribution function can be expressed as *probability mass function*. That is,
$$f_{\vec{X}}(x_1, x_2 \dots, x_n) = P(X_1 = x_1, X_2 = x_2 \dots, X_n = x_n)$$
 - If X_i s are continuous random variables, the joint distribution function can be defined as *probability density function* such that
$$P(X_1 \in (a_1, b_1), X_2 \in (a_2, b_2), \dots, X_n \in (a_n, b_n)) = \int_{a_n}^{b_n} \int_{a_{n-1}}^{b_{n-1}} \dots \int_{a_1}^{b_1} f_{\vec{X}}(x_1, \dots, x_n) dx_1 \dots dx_n.$$

for any real numbers a_i, b_i where $a_i < b_i$.
 - marginal distribution:
 - If X_i is discrete random variable, then we have
$$f_{X_i}(x_i) = P(X_i = x_i) = P(X_i = x_i \text{ and } X_j \in (-\infty, \infty) \text{ for all } j \neq i)$$

$$= \sum_{\text{all } x_1} \dots \sum_{\text{all } x_{i-1}} \sum_{\text{all } x_{i+1}} \dots \sum_{\text{all } x_n} f_{\vec{X}}(x_1, x_2, \dots, x_n).$$
 - If X_i is continuous random variable, we consider
$$F_{X_i}(x) = P(X_i \leq x) = P(X_i = x \text{ and } X_j \in (-\infty, \infty) \text{ for all } j \neq i)$$

$$= \int_{-\infty}^x \left(\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_{\vec{X}}(x_1, x_2, \dots, x_n) dx_1 \dots dx_{i-1} dx_{i+1} dx_n \right) dx_i$$

Then the probability density function of X_i can be computed as
$$f_{X_i}(x) = F'_{X_i}(x) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_{\vec{X}}(x_1, \dots, x_{i-1}, \textcolor{red}{x}, x_{i+1}, \dots, x_n) dx_1 \dots dx_{i-1} dx_{i+1} dx_n$$
 - change of variable formula

- g 是将 x 转变为 y 的函数，本质上是关于 x 的函数， h 是 g 的反函数， $y = Ax$ ，雅各比矩阵保证了将 x 转化为 y 的坐标系

$$f_Y(y_1, y_2, \dots, y_n) = f_X(h_1(y_1, \dots, y_n), \dots, h_n(y_1, \dots, y_n)) \left(\frac{1}{|J|} \right).$$

- Expected value = mean, covariance matrix = cov, 其实也就是之前的 Sigma
 - property:
 - $E[AX + b] = AE[X] + b$
 - $Cov(AX + b) = ACov(X)A^T$

- Multi-normal distribution

- 其实也就是 Gaussian vector

We let $\vec{Z} = \begin{pmatrix} Z_1 \\ \vdots \\ Z_n \end{pmatrix}$ be a vector of independent standard normal random variables (i.e. $Z_i \sim N(0, 1)$). We define

$$\vec{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} \mu_1 + a_{11}Z_1 + \dots + a_{1n}Z_n \\ \mu_2 + a_{21}Z_1 + \dots + a_{2n}Z_n \\ \vdots \\ \mu_n + a_{n1}Z_1 + \dots + a_{nn}Z_n \end{pmatrix} = \underbrace{\begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}}_{\vec{\mu}} + \underbrace{\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \ddots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}}_A \underbrace{\begin{pmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_n \end{pmatrix}}_{\vec{Z}}.$$

Then we say $\vec{Y} = \vec{\mu} + A\vec{Z}$ follows multivariate normal distribution with mean $\mathbb{E}[\vec{Y}] = \vec{\mu}$ and covariance matrix $\Sigma = Cov(\vec{Y}) = ACov(\vec{Z})A^T = AA^T$. We write $\vec{Y} \sim N_n(\vec{\mu}, \Sigma)$.

- 这里注意，多元正态除了 Gaussian vector 的定义还有一种是可以写成 $X = \mu + AZ$ 的形式， $AA^T = \Sigma$ ，其中 Z 是互相独立的，但如果不是 Gaussian vector，而仅仅只是两个任意的正态分布，那么是不能找到这样的 A 的，例子： X 正态， ϵX 正态（LC chapter4 的例子），不能找到这样的 A ，因为两者是相互独立的
- 同样，如果 (X, X) 按照 Gaussian vector 的定义确实是 Gaussian vector，那能不能写成那种形式呢，可以，按照后面给的公式就可以，一个自由变量，将它设为 0

- Properties

- if X 服从于多元正态分布，那么 $a_1X_1 + a_2X_2 + \dots + a_nX_n$ 服从正态分布
- if X 服从 $N_n(\mu, \Sigma)$, $Y = AX + b$ 服从 $N(A\mu + b, A\Sigma A^T)$
- if X 服从 $N_n(\mu, \Sigma)$, and Σ is invertible, $Y = (X - \mu)^T \Sigma^{-1} (X - \mu)$ 服从 χ^2 distribution with df = n
- Gaussian vector(大前提！！！): independent $\Leftrightarrow \text{cov} = 0$

We let $\vec{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_m \end{pmatrix}$ and $\vec{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$ be two vectors of random variables such that

$$\vec{W} = \begin{pmatrix} \vec{X} \\ \vec{Y} \end{pmatrix} \text{ follows multivariate normal distribution with } \vec{W} \sim N \left(\underbrace{\begin{pmatrix} \vec{\mu}_X \\ \vec{\mu}_Y \end{pmatrix}}_{\vec{\mu}}, \underbrace{\begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{XY} & \Sigma_{YY} \end{pmatrix}}_{\Sigma} \right)$$

Then \vec{X} and \vec{Y} are independent if and only if $\Sigma_{XY} = 0$ (zero matrix).

- 第四个性质并没有说 X 和 Y 必须独立，但是他们可以服从 multi normal distribution，但需要额外的条件，什么条件未知，如果我们由他们独立，那么他们就满足服从 multi normal distribution

We let \vec{X}_1 and \vec{X}_2 be two independent random variables such that $\vec{X}_1 \sim N_p(\vec{\mu}_1, \Sigma_{11})$ and $\vec{X}_2 \sim N_q(\vec{\mu}_2, \Sigma_{22})$. Then $\vec{X} = \begin{pmatrix} \vec{X}_1 \\ \vec{X}_2 \end{pmatrix}$ follows multivariate normal distribution with mean $\vec{\mu} = \begin{pmatrix} \vec{\mu}_1 \\ \vec{\mu}_2 \end{pmatrix}$ and covariance matrix $\begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$.

- 同样，如果独立，也可以获得两者相加也满足multi normal distribution, 但要注意他们的变量个数必须是相同的

- Central Limit Theorem

We let $\vec{X}_1, \vec{X}_2, \dots, \vec{X}_n$ (where $\vec{X}_i \in \mathbb{R}^p$) be n independent and identically distributed random variables with $\mathbb{E}[\vec{X}_i] = \vec{\mu}$ and $Cov(\vec{X}_i) = \Sigma$. Then when $n \rightarrow \infty$, $\frac{1}{n}(\vec{X}_1 + \vec{X}_2 + \dots + \vec{X}_n)$ is approximately multivariate normal with mean $\vec{\mu}$ and covariance matrix $\frac{1}{n}\Sigma$. $\sim N(\vec{\mu}, \frac{1}{n}\Sigma)$

- Drawbacks

- 如果我们选择利用多元正态分布来model stock return, 就可以得到每个stock都必须满足正态分布, 但这个事情是非常questionable的, 因为会有tail很大的stock, 也就不满足正态分布, 所以multi normal distribution is not a good fit

- Multi t distribution

- 定义的来源

- 我们在上面提到了, 并不是每个股票都满足normal distribution, 所以我们需要考虑另一种model, 在topic 1 我们已经学到了t distribution其实是蛮好的, 所以我们猜想能不能有multi t distribution 呢?
 - 然后我们直觉来看, 模仿multi normal distribution, 我们可以将multi t distribution 写成 $Y = \mu + AT_v$ 的形式, 但注意, T_v 很麻烦, 他每一个都有不同的自由度, 首先我们并不能确实这个的covariance是什么, 其次还记得我们之前求df吗, ad-hoc, it's so tedious!!!

$$\begin{pmatrix} T_{v1} \\ T_{v2} \\ \vdots \\ T_{vn} \end{pmatrix}$$

- 于是我们转化了一下multi t distribution 的形式, 其中 μ 称为location vector, v is the degree of freedom, 要求每个股票的自由度必须相同, 控制tail, v 越小, tail 越大, Z , 本来应该写成AZ, 为了方便, 我们现在将其搞成一个基础分布, 服从 normal distribution($0, \Lambda$), Z 之间independent

$$\vec{Y} = \vec{\mu} + \sqrt{\frac{v}{\chi^2_v}} \vec{Z},$$

- properties

- mean= μ , $cov(Y) = \frac{v}{v-2} \Lambda, v > 2$
- wY 也服从 t distribution: 这就对我们求portfolio return 有了帮助, 如果股票的return 服从 multi t distribution, 那么portfolio return 服从t 分布

- monte Carlo simulation

- 和single stock 的monte Carlo 的步骤类似, 这里我们要choose independent的 $X = (X_1, X_2, \dots, X_n)$, 如果以实例来看, 我们需要抽取n个股票的m个数据, 我们一直知道n个股票服从于multi t distribution ,而且他们之间很可能存在correlation, 但如果我们只是简单的抽取sample, 计算机是会independent抽取的, 这就和我们的假设不符, 所以我们要想一个办法独立抽取然后组成可能相关的multi t distribution
- 这就用到了之前的multi t $Y = \mu + \sqrt{\frac{v}{\chi^2_v}} Z$, Z 服从multi normal ,所以我们按照multi normal 的方法抽取Z, 然后抽取独立的 χ^2_v , 就可以获得Y了 (这对multi normal distribution也同样适用)
- process
 - multi normal distribution
 - 我们已知一个sample data, 假设X服从于multi normal ,我们可以用sample data 计算出应该服从的 μ, Σ , 然后我们可以重新抽取样本, $X = \mu + AZ, AA^T = \Sigma$, 可以计算出A, 然后只要抽取Z, 每个 Z_i 服从独立标准正态分布, 就可以计算了
 - multi t distribution
 - 和multi normal 一样, 我们已知 μ, Σ, v , 如何计算multi t 的A? 只要抽取一个随机的 χ^2_v , 就可以知道A, 然后抽取Z, Z 服从multi normal distribution, 抽取multi normal 我们在前面讲了, 只是这里的Z服从multi normal $(0, \Lambda)$, 然后可以计算 X_i , 再继续抽取多个 X_i , 就可以monte Carlo simulation

- Linear regression model

- 我们在上面确实推出了一个较为良好的model, 但是市场是在变化的, 所以model也应该变化, 而且每个model和每个model之间也应该是correlated的, 所以其实multi t 也并不是一个perfect 的model, 我们提出了linear regression model; 对于一个模型来说, 我们在构建的过程中会先猜想很多影响因素, 然后我们需要首先选择到底需要哪些变量, 真正有意义的变量有哪些? 有意义之后, 但是如果数量很多怎么办, 可不可以简化计算 (注意和减少变量数量不是一个概念! ! !) -->主成分分析
- 我们有学到, stock return 应该和很多因素有关系, eg. interest rate, inflation, GDP etc. 而且会有一些unexpected events, 于是我们建立了stock return 和多个因子之间的model, 为了简单, 我们将其定为线性模型;
 - response variable, 因变量
 - predictor variables, 自变量
 - ϵ_t , random error, 我们规定 $E[\epsilon_t] = 0, Var[\epsilon_t] = \sigma^2$, 为什么规定期望为0呢, 因为这种unexpected events有好有坏, 期望为0代表了我们对这种unexpected events have no information
 - β_i coefficients
- 有几种特殊的模型:
 - CAPM: r_i 是我们想要预测的stock return,

$$\overline{r_i} - r_f = \beta(\overline{f_M} - r_f) + \epsilon_i,$$

- 如何determine the coefficients?

- 我们选择让RSS最小, y_i 是真实值, $\hat{y}_i = \sum \beta_i y_i$ 是我们预测的值(不包含 ϵ)

$$RSS_p = \sum_{i=1}^n \left(\underbrace{y_i}_{\text{actual value}} - \underbrace{\left(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} \right)}_{\substack{\text{predicted value} \\ (\text{Linear regression})}} \right)^2$$

- 如何求解?

- 求导

To express the equations in a more convenient way, we define the following matrices:

$$\hat{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \in \mathbb{R}^{p+1}, \quad \vec{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n, \quad X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix} \in \mathbb{R}^{n \times (p+1)}.$$

Then the above system of equations can be expressed as

$$\begin{cases} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip}) = 0 \\ \sum_{i=1}^n x_{i1} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip}) = 0 \\ \vdots \\ \sum_{i=1}^n x_{ip} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip}) = 0 \end{cases} \Rightarrow \begin{cases} (1 \ 1 \ \dots \ 1) (\vec{Y} - X \hat{\beta}) = 0 \\ (x_{11} \ x_{21} \ \dots \ x_{n1}) (\vec{Y} - X \hat{\beta}) = 0 \\ \vdots \\ (x_{1p} \ x_{2p} \ \dots \ x_{np}) (\vec{Y} - X \hat{\beta}) = 0 \end{cases}$$

$$\Rightarrow X^T (\vec{Y} - X \hat{\beta}) = \vec{0} \Rightarrow \hat{\beta} = (X^T X)^{-1} X^T \vec{Y}.$$

- 但实际上, 仅仅只有一个point estimator 是非常有限的, point estimator 确实对我们目前的sample 的拟合in-sample是非常好的, 但对于out of sample的估计有可能并不好, 在这种情况下, 我们就需要寻找一个interval estimate, 这就涉及了他的四个性质

- 方程:

$$y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_p x_{1p} + \varepsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_p x_{2p} + \varepsilon_2$$

$$y_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_p x_{np} + \varepsilon_n$$

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} = n \times (p+1) \quad p+1 \leq n$$

$$\text{rank}(X^T X) = \text{rank}(p+1 \times p+1) = p+1 \quad \text{is invertible}$$

如果2个分子有3个X因素?
有多少解? 某个X可能完全用不到。

- 但在此之前, 我们对性质的推导过程进行了一些假设

- ϵ_i are independent, 服从于 $N(0, \sigma^2)$ --> mean=0 意味着 stay neutral, 为什么 independent 呢? 其实当然可以有 cov, 而且有 cov 的情况还很常见, 但是为了 simplify, 我们将其设为 independent 的
- X has full rank $p+1$
- property 1: $\hat{\beta} \sim N_{p+1}(\vec{\beta}, \sigma^2 (X^T X)^{-1})$

- proof

proof of 1: 如何证明从多元正态分布 $X = \mu + Az$

$$\hat{\beta} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (X\beta + e)$$

$$= \beta + \underbrace{(X^T X)^{-1} X^T e}_A.$$

$$AA^T = (X^T X)^{-1} X^T e \cdot e^T X (X^T X)^{-1} = e^T (X^T X)^{-1} = \text{cov}(\hat{\beta})$$

接下来我们对下面的推导进行些简化:

$$\hat{\beta} = X(X^T X)^{-1} X^T Y = H Y \quad \text{令 } H = X(X^T X)^{-1} X^T$$

property of H:

$$H^2 = H \quad H^2 = H^T \quad X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = H$$

$$H^T = H \quad H^T = X(X^T X)^{-1} X^T = H$$

$$\hat{\epsilon} = Y - \hat{\beta} = Y - H Y = (I - H)Y$$

- property 2: s^2 is the unbiased estimator of σ^2 , $s^2 = \frac{1}{n-p-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

- 注意这里的 s^2 不再是样本差了，因为分母不再是数目了而是 $n-p-1$

- proof

proof of 2:
target: $E[S^2] = \sigma^2$

由已知观察 $\sum (y_i - \hat{y}_i)^2 = \sum e_i^2$

$$E[\sum e_i^2] = \sum E[e_i^2] = \sum \text{Var}(e_i)$$

$$E[\sum e_i^2] = E[(y_i - \hat{y}_i)]^2 = E[X\beta + e_i - X\beta]^2 =$$

$$E[(e_i - (I-H)e_i)]^2 = \text{cov}(e_i) = \sigma^2 I$$

$$\therefore \text{cov}(e_i) = (I-H)^T e^2 I (I-H) = e^2 (I-H)$$

$$\therefore E[\sum e_i^2] = \text{tr}(\text{cov}(e_i)) = e^2 \text{tr}(I-H)$$

$$\text{tr} H = \text{tr}(X(X^T X)^{-1} X^T) = \text{tr}(X^T X(X^T X)^{-1}) = \text{tr}(I_p) = p$$

$$X(X^T X)^{-1} X^T = [n \times (p+1)] \times [n \times (p+1)] \times [(p+1) \times n] = [n \times (p+1)] \times [(p+1) \times n] = n \times n.$$

$$\therefore \text{tr}(H) = \text{tr}(AB) = \text{tr}(BA)$$

$$I_p = [(p+1) \times (p+1)] \times [(p+1) \times (p+1)] = (p+1) \times (p+1)$$

$$\therefore E[\sum e_i^2] = e^2 (n-p-1)$$

$$\therefore E[S^2] = E[\frac{\sum e_i^2}{n-p-1}] = \sigma^2 \quad (\text{s}^2 \text{是} e^2 \text{的无偏估计})$$

- property 3: $\hat{\beta}$ and s^2 is independent

- 首先肯定不可以直接求probability function, 太麻烦而且我们目前是求不出来的, 那怎么做呢, 我们可以找俩个等价的, 其中分别都只和其中一个相关
- $\hat{\beta}$ 和谁相关? X, Y, f ; s^2 和谁相关? e , 所以我们只需要证明 $\hat{\beta}$ and e independent
- 证明他们两个independent也不能用probability function
- proof

proof of 3. $\hat{\beta} \leftarrow X\hat{\beta} - X\beta$ $\hat{S} \leftarrow e_i$

$$\begin{aligned}\hat{e} &= (I - H)\varepsilon = (I - H)GZ. \\ X\hat{\beta} - X\beta &= HY - X\beta = HX\beta + HE - X\beta \\ &= X(X^T X)^{-1} X^T X\beta + HE - X\beta \\ &= HE = H\hat{e} Z\end{aligned}$$

Why \hat{e} independent? Gaussian vector + cov=0

$$\begin{pmatrix} \hat{e} \\ f \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \underbrace{\begin{pmatrix} \sigma(I-H) \\ \sigma H \end{pmatrix}}_{\text{cov}} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

$$\begin{aligned}\text{cov}(\hat{e}, f) &= \begin{pmatrix} \sigma^2(I-H) \\ \sigma^2 H \end{pmatrix} \cdot I_n \cdot \begin{pmatrix} \sigma(I-H) & \sigma H \end{pmatrix} \\ &= \sigma^2(I-H) + \sigma^2 H^2 \\ &\approx \underbrace{\sigma^2}_{\text{small}}.\end{aligned}$$

$\therefore \hat{e}$ independent.

- property 4: $\frac{RSS_p}{\sigma^2} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sigma^2} \sim \chi^2_{n-p-1}$
 - 如何证明服从卡方分布呢？证明是可以写成多个正态分布的平方求和的形式，有多少个正态分布求和这个卡方分布的自由度就是多少
 - proof

proof of 4 如何证明服从卡方分布？可以写成多个正态分布的平方和。

我们先找一个 $I-H$ 的特征值。

$$\begin{aligned}(I-H)v &= \lambda v \\ (I-H)(I-H)v &= \lambda(I-H)v \\ \lambda v &= \lambda^2 v \quad \lambda(v) = v \quad \therefore \lambda = 0 \text{ 或 } \lambda = 1\end{aligned}$$

因为 $v(I-H) = n-p-1$

\therefore 有 $n-p-1$ 个 0, $p+1$ 个 1

$$I-H \approx \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_{n-p} \end{pmatrix} \text{ 且 } \lambda_{n-p} = 0$$

且 $e^T e = (\alpha^T e)^T (\alpha^T e) = \alpha^T e^T \alpha^T e = \alpha^T D \alpha^T e = D^T \alpha^T e = D^T \alpha^T \alpha = \alpha^T \alpha$

$\alpha^T e = \alpha^T (I-H)z = \alpha^T D \alpha^T \alpha^T z = D^T \alpha^T \alpha^T z = D^T \alpha^T \alpha = \alpha^T \alpha$

$e^T e = (\alpha^T D^T \alpha^T z)^T (\alpha^T D^T \alpha^T z) = \alpha^T D^T \alpha^T z^T \alpha^T D^T \alpha^T z = \alpha^T \alpha$ 只要证 $\alpha^T \alpha$ 也是标准正态
independent

$E[\alpha^T \alpha] = E[\alpha^T \alpha] = 0$

$\text{Var}(\alpha^T \alpha) = \alpha^T \alpha - \text{Var}(\alpha^T \alpha) = 1$

步而 independent. $\hat{e} = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix} \xrightarrow{\text{独立}} \hat{e}^T \hat{e} \text{ 独立且同分布}$

从而证 $\text{cov}(\hat{e}^T \hat{e}) = 0$

$\text{cov}(\hat{e}^T \hat{e}) = \alpha^T \text{cov}(\hat{e}) \alpha = \alpha^T I_n \alpha = \alpha^T \alpha = n-p-1$

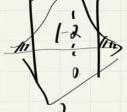
- 如何选择变量？
 - 如果变量太多，多一个变量就多一个需要估计的 β , 就会是 error 变大；但少一个变量就可能少了信息，所以如何选择一个最合适的变量呢？
 - 这个变量有没有意义 $\beta = 0$? -- hypothesis test: t-test/F-test
 - Method 1: T-test

- process

... < ↑ ⊞ ⊖ ⊛ ⊚ +

$$H_0: \beta_i = 0 \quad H_1: \beta_i \neq 0 \quad (\text{two-sided test})$$

$$\hat{\beta}_i \sim N_{p+1}(\hat{\beta}_i, C_{ii}^2(X^T X)^{-1}) \quad \begin{matrix} X: n \times p+1 \\ X^T: (p+1) \times n \end{matrix} \quad C = (X^T X)^{-1}$$

$$\hat{\beta}_i \sim N\left(\hat{\beta}_i, C_{ii}^2\right) \Rightarrow \frac{\hat{\beta}_i - \beta_i}{\sqrt{C_{ii}}} \sim N(0, 1) \Rightarrow \frac{\hat{\beta}_i - \beta_i}{\sqrt{C_{ii}}} \sim t_{n-p-1}$$


$$P(|T| \leq T^*) \geq 1 - \alpha$$

$$P(|T| \geq T^*) \leq \alpha. \quad T^* \text{ is critical value.}$$

① 看 $|T| > T^*$, reject H_0

② 计算 $P(|T| \leq T_0) \approx P \leq \alpha$. 说明 $|T_0| > T^*$, reject H_0

- Method 2: F-test

- 不再是一个一个检验了，我们之前不是有p个变量吗，我用p个variable做了一个 regression, 然后再用r个variable 做了一个regression, 如果可以简化成r个变量，那么这两套的方差应该是相近的，方差服从什么分布? χ^2 , 两个 χ^2 相比服从什么分布? F分布, 趋于1
- 真的吗? 当然是假的!!! 两个独立的 χ^2 分布相比才是F分布, 怎么让他们独立呢? 需要modify

... < ↑ ⊞ ⊖ ⊛ ⊚ +

$$\hat{y}_{it} = \beta_0 + \beta_1 x_{1it} + \beta_2 x_{2it} + \dots + \beta_p x_{pit}$$

$$\hat{y}'_{it} = \beta_0 + \beta_1 x_{1it} + \beta_2 x_{2it} + \dots + \beta_r x_{rit}$$

$$S_e^2 = \frac{1}{n-p-1} \sum_{i=1}^n (y_{it} - \hat{y}_{it})^2 \quad S_o^2 = \frac{\sum_{i=1}^n (y_{it} - \hat{y}'_{it})^2}{n-r-1} = \frac{RSS_r}{n-r-1} = \frac{C^2}{n-r-1}$$

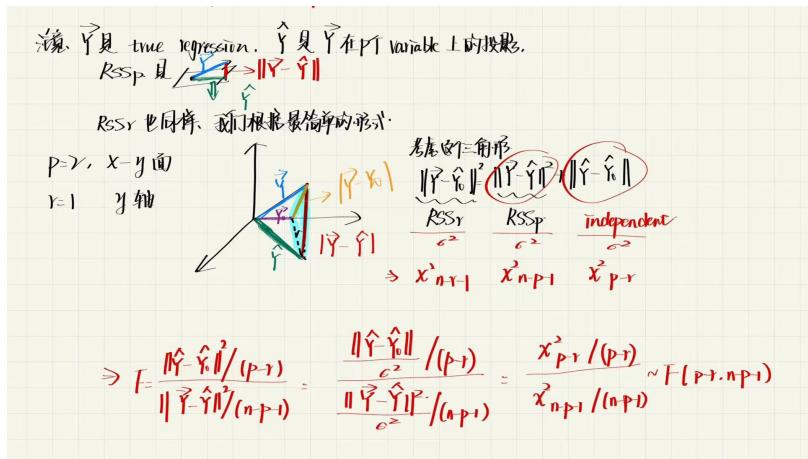
$$= \frac{RSS_p}{n-p-1} = \frac{RSS_p}{n-p-1} = \frac{C^2}{n-p-1} \quad \begin{matrix} \sim \chi^2_{n-p-1} \\ = z_1^2 + \dots + z_{n-p-1}^2 \end{matrix}$$

$$\sim \chi^2_{n-p-1} = z_1^2 + \dots + z_{n-p-1}^2 \quad \begin{matrix} \sim \chi^2_{n-r-1} \\ = z_1^2 + \dots + z_{n-r-1}^2 \end{matrix}$$

z_i 是 independent.

若转化为这种形式就具有 independent 的.

- modify 之后, 趋于0



- hypothesis

The null hypothesis and alternative hypothesis are set to be

$$\begin{cases} H_0: \beta_{r+1} = \beta_{r+2} = \dots = \beta_p = 0 & \text{Some of these factors are irrelevant} \\ H_1: \text{otherwise} \end{cases}$$

- H_0 的含义是，所有的都等于0， H_1 指只要有某些 $r+1$ 到 p 不等于0就可以，但我们不知道具体是哪些，所以其实还需要T-test
- 如果F足够大，则拒绝 H_0

- Method 3: Mallow's C_p statistics

- 其实是另一种方法来判断哪些变量的参数等于0，哪些不等于0，这个方法和T-test、F-test是独立的
- 我们做回归的目的是什么？predict，我们希望predict的error最小，所以如果我们用一部分数据做了回归之后，在用 X_{new} 做预测时，希望 $\min \sum \{y_{new} - \hat{y}_{new}\}^2$

$$\begin{aligned}
 & E \left[\sum_{t=1}^m (E[y_t^{new}] - \hat{y}_t^{new})^2 \right] \\
 & \approx \sum_{t=1}^m E \left[(E[y_t^{new}] - \hat{y}_t^{new})^2 \right] = E[y_t^{new}]^2 - E[y_t^{new}]E[\hat{y}_t^{new}] \\
 & = \sum_{t=1}^m E \left[(y_t^{new} - \hat{y}_t^{new})^2 \right] = PE \\
 & \text{其中 } P \leftarrow \sum_{t=1}^m E \left[\left(\sum_{j=0}^p \beta_j x_{tj} + \epsilon_t^{new} \right)^2 \right] \quad \text{我们假设 } x^{new} = x \\
 & \text{而已} \\
 & \text{不使用} \\
 & \text{而} \\
 & \text{已} \\
 & \text{不} \\
 & \text{使} \\
 & \text{用} \\
 & = \sum_{t=1}^m E \left[\left(\sum_{j=0}^p (\beta_j - \hat{\beta}_j) x_{tj} + \epsilon_t^{new} \right)^2 \right] \quad \text{如果新的事物也会发生} \\
 & = \sum_{t=1}^m E \left[\sum_{j=0}^p (\beta_j - \hat{\beta}_j) x_{tj}^2 \right] + E[\epsilon_t^{new}]^2 = \text{Var}(\epsilon_t^{new}) + (E[\epsilon_t^{new}])^2 \\
 & = \sum_{t=1}^m E \left[\sum_{j=0}^p (\beta_j - \hat{\beta}_j) x_{tj} \right]^2 + E[\epsilon_t^{new}]^2 = \sigma^2 \\
 & + E \left[\sum_{j=0}^p (\beta_j - \hat{\beta}_j) x_{tj} \cdot \epsilon_t^{new} \right] = 0 \quad \text{因为 } \epsilon_t^{new} \text{ 是 unexpected error, 而且新旧的数据预测时 have no info.} \\
 & = \sum_{t=1}^m E \left[\left(\sum_{j=0}^p (\beta_j - \hat{\beta}_j) x_{tj} \right)^2 \right] + n\sigma^2 \\
 & = \boxed{Error + n\sigma^2 = PE} \quad \text{说明不独立, 证明 model 错误.} \\
 & \text{即 } prediction \text{ error in the past} = prediction \text{ error in the future} \\
 & \text{我们设 } PE = E[RSP] = 2\sigma^2(p+1) \\
 & \therefore Error + n\sigma^2 - E[RSP] = 2\sigma^2 p + 2\sigma^2 \\
 & Error = 2\sigma^2 p + 2\sigma^2 + E[RSP] - n\sigma^2 = \sigma^2 E \left[2(p+1) - n + \frac{RSP}{\sigma^2} \right]
 \end{aligned}$$

- Remark:

- σ^2 未知，我们可以用 s_k^2 来代替
- 最好的model就是最小的 C_p statistics

- χ^2 分布的均值是 df, 方差是 $2 * \text{df}$, $E[RSS_p] = (n - p - 1) * \sigma^2$, $E[C_p] = p + 1$, 也就是说, 我们选择 C_p 最小的, 就可以解决 overfitting (变量冗余) 的情况
- 以上所有我们的结论都是基于一个最基本的假设来讲的: 这个 model 包含所有的 reliable 的变量, 如果有遗漏的变量, 那么上面的 $E[RSS_p] = (n - p - 1) * \sigma^2$, $E[C_p] = p + 1$ 并不成立, 而变成了下面的公式: (这里不证明)

$$\mathbb{E}[RSS_p] = \sum_{t=1}^n (\mathbb{E}[y_t - \hat{\beta}^T x_t])^2 + (n - p - 1)\sigma^2.$$

not true

- 我们的判断标准是最小的 C_p , 但实际上我们在这里做的只是计算 C_p , 你看 error 的公式其实是 C_p 的期望的最小值, 而不是我们真正用的仅仅只是选择最小的 C_p , 所以有可能我们得到的 model 再进行 T-test, 是可能会有 irrelevant 的 variable 存在的--> 保证了包含了所有的 relevant 的 variable, 但没有保证不包含所有 irrelevant 的 variable

- Lemma 的证明

The handwritten proof consists of two parts:

Top Part:

$$\begin{aligned} \mathbb{E}[RSS_p] &= \sum_{t=1}^n \mathbb{E}[(y_t - \hat{y}_t)^2] = \sum_{t=1}^n \mathbb{E}[(y_t - \hat{y}_t)^2] \\ &= \sum_{t=1}^n \mathbb{E}[y_t^2] - 2\mathbb{E}[y_t \hat{y}_t] + \mathbb{E}[\hat{y}_t^2] \\ &\quad - \mathbb{E}[y_t^2] + 2\mathbb{E}[y_t \hat{y}_t] - \mathbb{E}[\hat{y}_t^2] \end{aligned}$$

由 assumption: $X = X^{\text{new}} \Rightarrow \hat{y}_t = \hat{y}_t$

而我们知道 $y_t = X_t \beta + \varepsilon_t$. ε_t 和 $\varepsilon_t^{\text{new}}$ 服从于 $N(0, \sigma^2)$.

$\Rightarrow y_t$ 和 \hat{y}_t 服从于 same distribution

$\Rightarrow \mathbb{E}[y_t] = \mathbb{E}[\hat{y}_t]$

$\mathbb{E}[RSS_p] = \sum_{t=1}^n 2\mathbb{E}[y_t \hat{y}_t - y_t \hat{y}_t^{\text{new}}]$

$= \sum_{t=1}^n 2\mathbb{E}[y_t (\hat{y}_t - \hat{y}_t^{\text{new}})]$

$= \sum_{t=1}^n 2\mathbb{E}[y_t (\varepsilon_t - \varepsilon_t^{\text{new}})]$

$= 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T \hat{y}] = 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T X \hat{\beta}]$

$= 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T X (X^T X)^{-1} X^T \hat{y}] = 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T H (X \hat{\beta} + \varepsilon)]$

$= 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T H X \hat{\beta} + 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T H \varepsilon]]$

$= 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T H X \hat{\beta}] + 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T H \varepsilon]$

$= 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T H X \hat{\beta}] + 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T H \varepsilon]$

$= 2\mathbb{E}[(\varepsilon - \varepsilon^{\text{new}})^T H X \hat{\beta}] + 2\mathbb{E}[\varepsilon^T H \varepsilon] - \mathbb{E}[\varepsilon^T \varepsilon]$

$= 0$ $\because [n \times 1] \times [n \times n]$

$X (n \times (p+1)) \times (p+1) \times n$

$(n \times (p+1)) \times ((p+1) \times n) = (n \times (p+1)) \times (p+1) \times n = n \times n$

$\Rightarrow \mathbb{E}[H \varepsilon] = (1 \times n) \times (n \times n) \times (1 \times n)$

$= (1 \times n) \times (n \times n) = 1 \times 1 \therefore \mathbb{E}[H \varepsilon] \text{ just a number.}$

不能直接将 H 提出来, 但是做一个 $|x| \geq \text{number}$.

$\Rightarrow \mathbb{E}[H \varepsilon] = \text{tr}(\mathbb{E}[H \varepsilon])$

$\therefore \mathbb{E}[\varepsilon^T H \varepsilon] = \mathbb{E}[\text{tr}(\varepsilon^T H \varepsilon)] = \mathbb{E}[\text{tr}(\mathbb{E}[\varepsilon^T H \varepsilon])]$

$\mathbb{E}[\varepsilon^T H \varepsilon] = \mathbb{E}\left[\left(\begin{array}{c} \varepsilon_1^T H \varepsilon_1 \\ \vdots \\ \varepsilon_n^T H \varepsilon_n \end{array}\right)\right] = \mathbb{E}[\mathbb{E}[\varepsilon^T H \varepsilon]]$

同样 $\mathbb{E}[\varepsilon^T H \varepsilon] = \mathbb{E}[\text{tr}(\varepsilon^T H \varepsilon)] = \mathbb{E}[\text{tr}(\mathbb{E}[\varepsilon^T H \varepsilon])]$

$$\begin{aligned}
 & (\text{Im}(\Phi)) \cup (\text{pr}_1) \times (\text{pr}_1) = (\text{Im}(\Phi)) \vee (\text{pr}_1 \times \text{pr}_1) \\
 \Rightarrow & \text{E}^1 H_2 = (\text{Im}(\Phi)) \times (\text{Im}(\Phi)) \times (\text{Im}(\Phi)) \\
 & = (\text{Im}(\Phi)) \times (\text{Im}(\Phi)) = \text{Im}(\Phi) \quad \because \text{E}^1 H_2 \text{ just a number.} \\
 & \text{不能自己将自己提出来. 但要是 } \text{Im}(\Phi) \Rightarrow \text{number.} \\
 \Rightarrow & \text{E}^1 H_2 = \text{tr}(\text{E}^1 H_2) \\
 \therefore & \text{tr}[\text{E}^1 H_2] = \text{tr}[\text{tr}(\text{E}^1 H_2)] = \text{tr}[\text{tr}(\text{E}^1 H)] \\
 \text{tr}(\text{E}^1 H) & = \text{tr}\left[\left(\begin{array}{cccc} E_1 & E_{12} & \dots & E_{1n} \\ 0 & E_2 & & \\ \vdots & & \ddots & \\ 0 & 0 & \dots & E_n \end{array}\right) H\right] = E^1 \text{tr}(H) \\
 \text{同样 } \text{tr}[\text{E}^1 \text{new } H_2] & = \text{tr}[\text{tr}(\text{E}^1 \text{new } H_2)] = \text{tr}[\text{tr}(\text{E}^1 \text{new } E H)] \\
 \text{tr}(\text{E}^1 \text{new } E H) & = \text{tr}\left[\left(\begin{array}{cccc} E_1^{\text{new}} E_1 & & & \\ \vdots & \ddots & & \\ 0 & 0 & \ddots & \\ 0 & 0 & \dots & E_n^{\text{new}} E_n \end{array}\right) H\right] = 0. \\
 \therefore \text{tr}[\text{E}^1 H] & = 2E^2 \text{tr}(H) = 2E^2 \text{tr}[X(X^T X)^{-1} X^T] = 2E^2 \text{tr}(X^T X (X^T X)^{-1}) = 2E^2 (\text{pr}_1). \\
 X^T X & = (\text{pr}_1 \times \text{Im}(\Phi)) \times (\text{Im}(\Phi) \times \text{pr}_1) = (\text{pr}_1) \times (\text{pr}_1). \\
 & = \text{PE} = \text{tr}(\text{RSp}) \\
 \text{Lemma: 证明完毕.}
 \end{aligned}$$

- Method4: AIC

- C_p 实际上只考虑了expectation, AIC和 C_p 不同的点在于它考虑的是distribution, 而不是简单的expectation, $|g_Y(y) - f_Y(y)|$, 但这个公式是non-feasible的, 所以我们选择了average distance

$$D = \int_{-\infty}^{\infty} [ln f_Y(y) - ln g_Y(y)] f_Y(y) dy = E[ln f_Y(y) - ln g_Y(y)]$$

- Remark: 为什么取了ln?

- 可以simplify calculation
 - guarantee positive

$$\begin{aligned}
 D &= \int_{\text{do}}^{\text{tp}} [\ln f_Y(y) - \ln g_Y(y)] f_Y(y) dy = E[\ln f_Y(y) - \ln g_Y(y)] \\
 -D &= E[\ln g_Y(y)] - E[\ln f_Y(y)] = E\left[\ln \frac{g_Y(y)}{f_Y(y)}\right] \leq E\left[\frac{g_Y(y)}{f_Y(y)} - 1\right] \\
 &= \int \left[\frac{g_Y(y)}{f_Y(y)} - 1 \right] f_Y(y) dy = \int g_Y(y) - f_Y(y) dy = 1 - 1 = 0 \\
 \therefore -D &\leq 0, \quad D \geq 0
 \end{aligned}$$

- 这两条也就说明了这个D是个reasonable indicator
 - 我们希望D min, 而 $E[\ln f_Y(y)]$ 是一定的, 不会发生改变, 所以希望g max, g是什么? 是我们估计的y的distribution, 让distribution的期望最大又是什么? --> MLE
 - $E[\ln g_Y(y)]$ 这个很熟悉吧, 蒙特卡洛模拟!

$$E[\ln g_Y(y)] \approx \frac{1}{n} \sum_{i=1}^n \ln g_Y(y_i | \theta) = \frac{1}{n} \underbrace{\ln g_Y(y_1) g_Y(y_2) \dots g_Y(y_n)}_{= \frac{1}{n} \ell(\theta | y_1, \dots, y_n)}$$

- 经过证明，但实际上，这并不是一个无偏估计，所以经过证明我们能得到，这里的 k 是参数的数量，也就是 $p+2, p+1$ 个 β +一个 σ

$$\mathbb{E}\left[\frac{1}{n}l(\hat{\theta}|y_1, \dots, y_n)\right] \approx \mathbb{E}[\ln g_Y(y|\vec{\theta})]_{\vec{\theta}=\hat{\theta}} + \frac{k}{n}$$

$$\Rightarrow \mathbb{E}\left[\frac{1}{n}l(\hat{\theta}|y_1, \dots, y_n) - \frac{k}{n}\right] = \mathbb{E}[\ln g_Y(y|\vec{\theta})]_{\vec{\theta}=\hat{\theta}}.$$

- $\max g \rightarrow \max E[\frac{1}{n}l(\theta) - \frac{k}{n}] \rightarrow \max E[-\frac{1}{n}(-l(\theta) + k)] \rightarrow \min(-2l(\theta) + 2k) \rightarrow \min AIC = \min(-2l(\theta) + 2k)$
- 接下来我们计算着 $p+2$ 个参数的理论值：， β 就是我们之前的理论值， σ 通过求导取最值

$$L(\beta_0, \dots, \beta_p, \sigma^2) = \prod_{i=1}^n \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2}{2\sigma^2}}}_{g_{Yt}(y_i|\beta_0, \dots, \beta_p, \sigma^2)} = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} e^{-\sum_{i=1}^n \frac{(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2}{2\sigma^2}}$$

$$\Rightarrow l(\beta_0, \dots, \beta_p, \sigma^2) = \ln L(\cdot) = -\frac{n}{2} \ln(2\pi) - n \ln \sigma - \sum_{i=1}^n \frac{(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2}{2\sigma^2}.$$

$$\hat{\beta} = (X^T X)^{-1} X^T \vec{Y}.$$

$$\frac{\partial l}{\partial \sigma} = -\frac{n}{\sigma} \frac{\sum_{i=1}^n (y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{ij})^2}{\sigma^3} = 0 \Rightarrow \hat{\sigma} = \sqrt{\frac{RSS_p}{n}}.$$

$$l(\hat{\theta}) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \frac{RSS_p}{n} - \frac{RSS_p}{2 \left(\frac{RSS_p}{n} \right)} = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} - \frac{n}{2} \ln \frac{RSS_p}{n}.$$

$$AIC = -2 \left(-\frac{n}{2} \ln(2\pi) - \frac{n}{2} - \frac{n}{2} \ln \frac{RSS_p}{n} \right) + 2(p+2)$$

$$= n \ln(2\pi) + n + n \ln \frac{RSS_p}{n} + 2(p+2)$$

$$= n \ln(2\pi) + n + n \left(\ln \frac{RSS_p}{n} + \frac{2(p+2)}{n} \right) \sim \ln \frac{RSS_p}{n} + \frac{2(p+1)}{n}$$

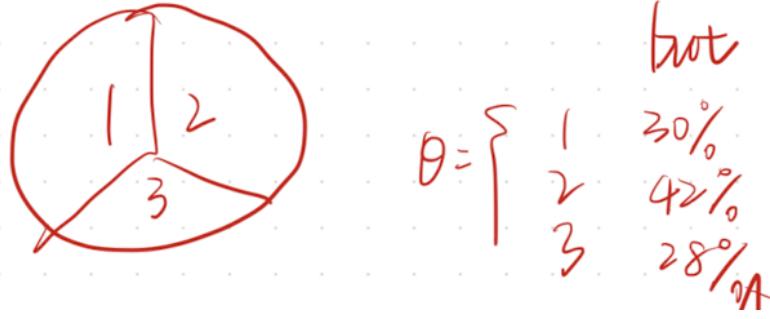
6^2必须一直有
78(p+1)可能会变

- 这里为什么把 $p+2$ 写成 $p+1$ ， 是因为 σ 是一直都有的，但其他的 β 可能会存在 irrelevant 的 factor，所以不一定都存在， \rightarrow 只有 $p+1$ 个 free parameters

- Method 5 : BIC

- BIC 和 AIC 的不同：

- 还记得吗，我们在 AIC 中应用了蒙特卡洛模拟，应用蒙特卡洛模拟的前提是必须有一个 large sample，如果没有这个前提，比如 IPO 数据，结果就不再 reliable
- 之前在数据很少的情况下我们怎么做的？ find some similar stocks to indicate something about this stock \rightarrow 但注意，这个 $\theta_i \neq \theta$ ，我们 treat θ as a random variable \rightarrow prior distribution，例如



- 一旦有了新的信息，我们就update这个distribution --> information updating
- process:
 - 我们需要一个prior distribution, 这个distribution感觉包含了所有的结果，且之后所有的结果可能不会再增加新的事件 --> $\pi(\theta) = 1(40\%), 2(30\%), 3(30\%)$, $f_{X^*}(x) = \int_{-\infty}^{+\infty} f_X(x|\theta)\pi(\theta)d\theta$, 这个函数是概率密度函数，也就是 $X^* = x$ 的概率是这个， x 是提前已知的，是我们的输入量，所以是对 θ 进行积分
 - 之后，经过了一段时间之后，我们可以获得一些新信息，这时就要更新我们的参数的分布，进而更新我们的新的概率密度函数，此时计算概率密度函数的公式不变，只是 $\pi(\theta)$ 变了

$$\pi(\theta|X_1, X_2, \dots, X_n) = \frac{\pi(\theta)f_X(X_1, X_2, \dots, X_n|\theta)}{\int_{-\infty}^{\infty} \pi(\theta)f_X(X_1, X_2, \dots, X_n|\theta)d\theta}.$$

- 我们怎么利用这个process来进行model selection?
 - 假设我们可以假设m个model，我们suppose每个model（其实每个model就是不同的parameter组合），也就是每个 θ 的初始概率是 $\frac{1}{m}$
 - 这里我们更改一下denotation, $\pi(\theta) = g_i(\beta_i)$
 - BIC score:

- 1

prior: $f_{\theta}(y_1, \dots, y_n | M_i) = \int_{\theta_i} \tilde{f}_{\theta}(y_1, \dots, y_n | \theta_i, M_i) g_i(\theta_i) d\theta_i$
 我们令 $h(\beta_i) = \ln(f_{\theta}(y_1, \dots, y_n | \beta_i))$ 且我们令 $\hat{H}(\beta_i)$ 表示 $h(\beta_i)$ 的二阶导数
 我们利用 Taylor expansion 对 $h(\beta_i)$ 在 β_i^* 处展开。

$$h(\beta_i) = h(\beta_i^*) + \frac{1}{2} (\beta_i - \beta_i^*)^T \left[\frac{\partial^2 h(\beta_i^*)}{\partial \beta_i \partial \beta_j} \right] (\beta_i - \beta_i^*) + \frac{1}{2} \left[\frac{\partial^2 h(\beta_i^*)}{\partial \beta_i \partial \beta_k} \right] (\beta_i - \beta_i^*)(\beta_k - \beta_k^*) \dots$$

$$\text{且 } \max_{\beta_i} f_{\theta} \Rightarrow \frac{\partial h(\beta_i)}{\partial \beta_i} = 0$$

$$h(\beta_i) \approx h(\beta_i^*) + \frac{1}{2} (\beta_i - \beta_i^*)^T H(\beta_i - \beta_i^*) \quad H(\beta_i^*) = \left(\frac{\partial^2 h(\beta_i^*)}{\partial \beta_i \partial \beta_k} \right) \text{ 称为 Hessian matrix}$$

$$\text{且 } H(\beta_i) \text{ 是半正定的。} \quad \nabla H(\beta_i^*) \text{ 且 } H(\beta_i^*) \text{ 一般不是负定的。}$$

$$\Rightarrow H(\beta_i^*) = H(\beta_i^*)^T \quad H(\beta_i^*) \text{ 是 positive definite.}$$

$$f_{\theta}(y_1, \dots, y_n | M_i) = \int_{\theta_i} e^{\ln(f_{\theta}(y_1, \dots, y_n | \theta_i, M_i))} d\theta_i = \int_{\theta_i} e^{h(\theta_i) - \frac{1}{2} (\theta_i - \theta_i^*)^T H(\theta_i^*) (\theta_i - \theta_i^*)} d\theta_i$$

Because Hessian matrix H / \hat{H} is symmetric. $\Rightarrow \hat{H}(\beta_i^*) = UDU^T$ semipositive

$$D = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_m \end{pmatrix} \lambda_i > 0. \quad U \text{ is orthogonal matrix. } UU^T = I. \quad \det(UU^T) = \det(U)\det(U^T) = 1$$

$$\therefore f_{\theta}(y_1, \dots, y_n | M_i) = e^{\ln(f_{\theta}(y_1, \dots, y_n | M_i))} \int_{\theta_i} e^{-\frac{1}{2} (\theta_i - \theta_i^*)^T UDU^T (\theta_i - \theta_i^*)} d\theta_i \quad \det U = 1$$

- 2

$$\begin{aligned}
& \text{let } y = \mathbf{y}^T(\beta_1 \ \beta_2) \\
&= e^{h(\beta_1)} \int_{\beta_1}^{\infty} e^{-\frac{1}{2} y^T \mathbf{y}} \det(\mathbf{y}) dy = e^{h(\beta_1)} \int_{\beta_1}^{\infty} e^{-\frac{1}{2} y^T \mathbf{y}} dy \\
&= e^{h(\beta_1)} \int_{\beta_1}^{\infty} e^{-\frac{1}{2} \sum_{i=1}^n y_i^2} dy = e^{h(\beta_1)} \prod_{i=1}^n \int_{\beta_1}^{\infty} e^{-\frac{1}{2} y_i^2} dy \\
&= e^{h(\beta_1)} g_i(\beta_1)^n \cdot \frac{1}{\det(\mathbf{y})} \cdot \ln \int_{\beta_1}^{\infty} e^{-\frac{1}{2} y^T \mathbf{y}} dy = \frac{1}{\det(\mathbf{y})} \ln \int_{\beta_1}^{\infty} e^{-\frac{1}{2} y^T \mathbf{y}} dy \\
&\therefore \ln f(y_1, \dots, y_n | \mathbf{m}) = \ln L(\beta_1) + \ln g_i(\beta_1) + \frac{1}{2} \ln(n) - \frac{1}{2} \ln(\det(\mathbf{y}))
\end{aligned}$$

constant $\ln \int_{\beta_1}^{\infty} e^{-\frac{1}{2} y^T \mathbf{y}} dy = m$

$$[\ln f]_{\text{nm}} = \frac{\partial^2 h}{\partial \beta_{\text{nm}} \partial \beta_{\text{nm}}} = \frac{\partial^2}{\partial \beta_{\text{nm}} \partial \beta_{\text{nm}}} [\ln L(\beta_1) + \ln g_i(\beta_1)] \text{ constant. } \frac{\partial^2}{\partial \beta_{\text{nm}} \partial \beta_{\text{nm}}} = \frac{\partial^2}{\partial \beta_{\text{nm}} \partial \beta_{\text{nm}}} \ln L(\beta_1)$$

$$\therefore [\ln f]_{\text{nm}} = -\frac{\partial^2}{\partial \beta_{\text{nm}} \partial \beta_{\text{nm}}} \ln L(\beta_1) = -\frac{\partial^2}{\partial \beta_{\text{nm}} \partial \beta_{\text{nm}}} \sum_{i=1}^n \ln f(y_i | \beta_1, \mathbf{m})$$

$$\ln L(\beta_1) = \frac{1}{n} \sum_{i=1}^n \ln f(y_i | \beta_1, \mathbf{m}) \quad \downarrow E[n \ln f(Y | \beta_1, \mathbf{m})]$$

$$\downarrow \frac{n^2}{2 \beta_{\text{nm}}} \text{ if } E[\ln f(Y | \beta_1, \mathbf{m})] \text{ is a fixed number.}$$

$$\det(\mathbf{y}) = n^{\frac{1}{2}} \det(\mathbf{y}^T). \quad [\ln f]_{\text{nm}} = -\frac{\partial^2}{\partial \beta_{\text{nm}} \partial \beta_{\text{nm}}} E[\ln f(Y | \beta_1, \mathbf{m})]$$

$$\therefore \ln f(y_1, \dots, y_n | \mathbf{m}) = \ln L(\beta_1) + \ln g_i(\beta_1) + \frac{1}{2} \ln(n) - \frac{1}{2} \ln(\det(\mathbf{y}))$$

• 3

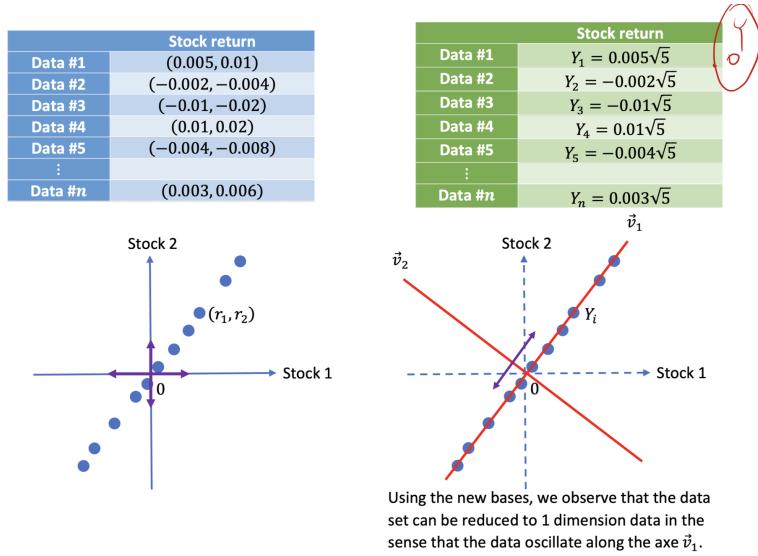
$$\begin{aligned}
& \ln f(y_1, \dots, y_n | \mathbf{m}) = \ln L(\beta_1) + \ln g_i(\beta_1) + \frac{1}{2} \ln(n) - \frac{1}{2} \ln(\det(\mathbf{y})) \\
& \text{因为 } n \text{ 是样本数量 } \Rightarrow n \text{ 是固定数.} \quad \text{根据.} \quad \text{根据.} \\
& \quad \text{根据.} \quad \text{根据.} \\
& \approx \ln L(\beta_1 | y_1, \dots, y_n) - \frac{1}{2} \ln n - \underbrace{\frac{1}{2} (\ln n - 2 \ln(\det(\mathbf{y})))}_{\text{BIC.}} \\
& \max \ln f(y_1, \dots, y_n | \mathbf{m}) \Rightarrow \min \text{ BIC.} \quad \text{AIC} \text{ 是 min.}
\end{aligned}$$

• Model selection

- Method 1: Forward selection
 - 刚开始考虑null model, 然后一个一个加variable
 - 刚开始先加p最小的, 即拒绝H_0 ($\beta \neq 0$) 的, 也就是F越大
- Method 2: Backward Selection
 - 刚开始考虑full model, 然后一个一个减variable
 - 刚开始减去 $\beta = 0$, 也就是不能拒绝H_0的, 即p最大的, F最小的
- Principal Component Analysis
 - reduce the dimension of the data
 - 比如这个数据, stock 1 和 stock 2 的每组数据都是有关系的, 都是1:2的关系, 所以我们只需要每组数据之间的系数

	Stock 1	Stock 2
Data #1	0.005	0.01
Data #2	-0.002	-0.004
Data #3	-0.01	-0.02
Data #4	0.01	0.02
Data #5	-0.004	-0.008
:	:	:
Data #n	0.003	0.006

- result



- 但问题是，实际上并没有这样perfect的dataset，如果我们reduce the dimension,就一定会产生error,所以我们reduce 的target就是使得误差最小
- assumption: 我们需要是的 $E[X_i] = 0$, 也就是每个dataset的期望必须是0, 如果期望不是0的话，他在坐标系上的中心点就不是坐标原点，我们是没有办法仅仅通过旋转坐标轴来使得所有的数据在更小的维度上
- process
 - 1

◀ ▶ ⌂

$X_i = \begin{pmatrix} X_{i1} \\ X_{i2} \\ \vdots \\ X_{ip} \end{pmatrix}$ CP. X_{ij} 是第 j 个 dataset 中第 i 个 stock 的 return. \vec{x}_i 是 x_i 的 $\vec{\cdot}$

举例: $X_i = X_{i1} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} + X_{i2} \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \dots + X_{ip} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$

我们希望找到多组系数 v_1, \dots, v_p (且向量不进行 reduce)

$$X_i = Y_{i1}\vec{v}_1 + Y_{i2}\vec{v}_2 + \dots + Y_{ip}\vec{v}_p \quad \|v_i\| = 1 \quad \text{target} \in \min \text{ error}$$

$$Y_{ij} = \underbrace{X_i^T}_{V_i} \vec{v}_j = V_i^T \vec{v}_j \quad \tilde{X}_i = Y_{i1}\vec{v}_1 + Y_{i2}\vec{v}_2 + \dots + Y_{im}\vec{v}_m, \text{ kmp}$$

$$E = \frac{1}{n} \sum_{i=1}^n \|X_i - \tilde{X}_i\|^2 = \frac{1}{n} \sum_{i=1}^n \|X_i - \left(\sum_{j=1}^m Y_{ij} v_j \right)\|^2.$$

$$\|X_i\|^2 = \|X_i - \tilde{X}_i\|^2 + \|\tilde{X}_i\|^2 \quad \text{max.}$$

constant: $\max_{i=1}^n \|X_i\|^2 = \max_{i=1}^n \sum_{j=1}^m Y_{ij}^2 = \max_{V_1, \dots, V_m} \sum_{i=1}^n \sum_{j=1}^m (V_j^T X_i)^2$

target: $\max_{V_1, \dots, V_m} \sum_{i=1}^n \sum_{j=1}^m (V_j^T X_i)^2 = \max_{V_1, \dots, V_m} \sum_{i=1}^n \sum_{j=1}^m (V_j^T X_i)(V_i^T X_i) = \max_{V_1, \dots, V_m} \sum_{i=1}^n \sum_{j=1}^m V_j^T (X_i X_i^T) V_i$

我们想利用 n 个线性性质, snap the sum

$$= \max_{V_1, \dots, V_m} \sum_{j=1}^m V_j^T \begin{pmatrix} X_{i1} X_{i1} & X_{i1} X_{i2} & \dots & X_{ii} X_{ip} \\ \vdots & \ddots & \ddots & \vdots \\ X_{ip} X_{ip} \end{pmatrix} V_j$$

$$\approx n \max_{V_1, \dots, V_m} \sum_{j=1}^m V_j^T [E(X_i X_i^T)] V_j \quad \text{kmp} \quad n \max_{V_1, \dots, V_m} \sum_{j=1}^m V_j^T \text{cov}(X) V_j$$

$$= n \max_{V_1, \dots, V_m} \text{Var}(V_j^T X) \quad \text{how to max.} \quad \text{Var}(V_j^T X)$$

• 2

我们从一个简单的 case 开始, 假设 V_i 到只有 1 dimension.

\Rightarrow Let $a = \text{cov}(X)$. $\Rightarrow \max_{V_1, \dots, V_i} V_i^T a V_i \quad \|V_i\|^2 = 1$

证明我们用 Lagrange Method definition:

The optimal solution (X_1^*, \dots, X_n^*) can be found by considering $f(x_1, \dots, x_n) - \lambda_1 g_1(\cdot) - \lambda_2 g_2(\cdot) - \dots - \lambda_m g_m(\cdot)$.

aim to solve $\min / \max f(x_1, \dots, x_n)$

subject to $\begin{cases} g_1(x_1, \dots, x_n) = 0 \\ g_2(x_1, \dots, x_n) = 0 \\ \vdots \\ g_m(x_1, \dots, x_n) = 0 \end{cases} \quad \|V_i\|^2 = 1 \quad \|V_i\|^2 - 1 = 0$

$\Rightarrow V_i^T V_i - \lambda_1 (\|V_i\|^2 - 1) = \sum_{j=1}^p \sum_{i=1}^n V_{ij} V_{ij} \lambda_j - \lambda_1 (\sum_{i=1}^n V_{ii}^2 - 1)$

由 first order condition:

$$\frac{\partial L}{\partial V_i} = \begin{pmatrix} \frac{\partial L}{\partial V_1} \\ \vdots \\ \frac{\partial L}{\partial V_p} \end{pmatrix} = 0 \Rightarrow \begin{cases} V_{11}^T V_{11} \lambda_1 + V_{12}^T V_{12} \lambda_2 + \dots + V_{1p}^T V_{1p} \lambda_p \\ V_{21}^T V_{21} \lambda_1 + V_{22}^T V_{22} \lambda_2 + \dots + V_{2p}^T V_{2p} \lambda_p \\ \vdots \\ V_{p1}^T V_{p1} \lambda_1 + V_{p2}^T V_{p2} \lambda_2 + \dots + V_{pp}^T V_{pp} \lambda_p \end{cases}$$

$$\Rightarrow \begin{cases} V_{11}^T V_{11} \lambda_1 + V_{12}^T V_{12} \lambda_2 + \dots + V_{1p}^T V_{1p} \lambda_p \\ V_{21}^T V_{21} \lambda_1 + V_{22}^T V_{22} \lambda_2 + \dots + V_{2p}^T V_{2p} \lambda_p \\ \vdots \\ V_{p1}^T V_{p1} \lambda_1 + V_{p2}^T V_{p2} \lambda_2 + \dots + V_{pp}^T V_{pp} \lambda_p \end{cases} = 0$$

$$\Rightarrow \begin{cases} V_{11}^T V_{11} \lambda_1 + V_{12}^T V_{12} \lambda_2 + \dots + V_{1p}^T V_{1p} \lambda_p \\ V_{21}^T V_{21} \lambda_1 + V_{22}^T V_{22} \lambda_2 + \dots + V_{2p}^T V_{2p} \lambda_p \\ \vdots \\ V_{p1}^T V_{p1} \lambda_1 + V_{p2}^T V_{p2} \lambda_2 + \dots + V_{pp}^T V_{pp} \lambda_p \end{cases} = 0$$

$$\therefore \boxed{\lambda_1 = \lambda_2 = \dots = \lambda_p}$$

∴ 如果 $\max \rightarrow$ 取最大的特征值 $\leq \max_{V_i} V_i^T a V_i = V_i^T \lambda_i V_i = \lambda_i V_i^T V_i = \lambda_i$

• 3

how to generalize? 其实应该预测 特征值.

$X_i = Y_{i1}\vec{v}_1 + \dots + Y_{im}\vec{v}_m$

$\max_{V_1, \dots, V_m} \sum_{i=1}^m V_i^T a V_i \quad \text{subject to } \sum_{i=1}^m \|V_i\|^2 = \|V\|^2 = 1$

这个真的非常 tedious. 因为一个具有非常大的数据集, 但是 $C^n \rightarrow$ try to simplify.

若 V_i 是 symmetric, \Rightarrow 特征值 orthogonal. $\langle V_i, V_j \rangle = 0$

\Rightarrow 这个 constraints 可以直接被 cancel out.

$\Rightarrow 1. \sum_{i=1}^m V_i^T a V_i - \sum_{i=1}^m \lambda_i \delta_{ii} (\|V\|^2 - 1)$

$\frac{\partial L}{\partial V_i} = 0 \Rightarrow \lambda_i V_i = \lambda_i V_i \Rightarrow \sum_{i=1}^m V_i^T a V_i = \sum_{i=1}^m V_i^T (\lambda_i V_i) = \sum_{i=1}^m \lambda_i$

• Remark:

- Y_{ij} principal components
- v_j factor loading
- Property

- $Var(Y_{ij}) = \lambda_j$
- $Cov(Y_{ij}, Y_{ik}) = 0$
- proof

Proof of the properties

We let $U = (\vec{v}_1, \vec{v}_2, \dots, \vec{v}_p)$ be a $p \times p$ matrix which the columns are the eigenvectors of Ω . Since Ω is symmetric, so \vec{v}_i s are orthogonal and $UU^T = I_p$.

Then the data \vec{X}_i (with no approximation) can be expressed as

$$\vec{X}_i = Y_{i1}\vec{v}_1 + Y_{i2}\vec{v}_2 + \dots + Y_{ip}\vec{v}_p = U\vec{Y}_i.$$

This implies that $\vec{Y}_i = U^T\vec{X}_i$.

Using the properties of covariance, we deduce that

$$Cov(\vec{Y}_i) = U^T \underbrace{Cov(\vec{X}_i)}_{=\Omega} (U^T)^T = U^T \left(U \underbrace{\begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_p \end{pmatrix}}_{=D} U^T \right) U = \underbrace{\begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_p \end{pmatrix}}_{=D}$$

This implies that $Var(Y_i) = Cov(Y_i, Y_i) = \lambda_i$ and $Cov(Y_i, Y_j) = 0$.

- total variance of $X_i = \sum_{i=1}^p \lambda_i = \text{tr}(D)$
- total variance of $\hat{X}_l = \sum_{i=1}^m \lambda_i = \text{tr}(D)$
- parallel shift component: affect in the same direction
- tilt component: short term is opposite with the long term
- curvature: mid -term is different from else