

# QUANTITATIVE FINANCE IN PRACTICE

MOHSIN JAVED

## 1. INTRODUCTION

Mathematical finance, as applied to the derivatives market in particular.

## 2. BOOKS, ARTICLES AND TOPICS WE SHOULD KNOW

- Black Scholes Model (Read Hull or Shreve)
- Change of measure
- Change of numeraire
- Forward measure
- Read Chapter 2 of Brigo/Mercurio
- For Quanto/Compo Read paper by Derman on Foreign indices
- How is a PDE equivalent to an expectation? (Feynman-Kac and all that)
- Forward and backward Kolmogorov equations
- Local volatility Model (Dupire equation, in particular)
- How the denominator and numerator of the local vol formula imply arbitrage conditions.
- Stochastic volatility model.
- conditional variance
- *Price is the cost of hedging*
- Understand the following two:
  - Analogy 1: Volatility – Theta – Gamma
  - Analogy 2: Correlation – Theta – Cross Gamma
- Understand the financial implications of the following and how much of a difference they drive from the simple Black Scholes model.
  - Funding Spreads
  - Dividends
  - Interest Rates
- Funding Valuation Adjustment (FVA), read paper by Piterbarg

---

*Date:* September 29, 2017.

### 3. RISK, RETURN AND SHARPE RATIO

Excess return, (i.e. actual return minus the risk free return) is proportional to risk (i.e. volatility).

Let  $r$  be the risk free return,  $\mu$  and  $\mu'$  be the actual returns of two portfolios with volatilities  $\sigma$  and  $\sigma'$ . It can be shown that the Sharpe ratio of each portfolio is the same [2]. That is,

$$(3.1) \quad \frac{\mu' - r}{\sigma'} = \frac{\mu - r}{\sigma} := \lambda$$

### 4. BLACK SCHOLES EQUATION WITH REPLICATION

At time  $t$ , a portfolio  $\Pi$  is formed by selling a call option for price  $C$ , which is delta hedged by buying  $\Delta$  of the underlying stock, at a cost of  $\Delta S$ :

$$\Pi = C - \Delta S.$$

At time  $t + dt$ , we have by Ito's formula:

$$\begin{aligned} d\Pi &= dC - \Delta dS \\ r\Pi dt &= \frac{\partial C}{\partial t} dt + \frac{\partial C}{\partial S} dS + \frac{1}{2} \frac{\partial^2 C}{\partial S^2} dS^2 - \Delta dS \\ r(C - \Delta S)dt &= \frac{\partial C}{\partial t} dt + \left( \frac{\partial C}{\partial S} - \Delta \right) dS + \frac{1}{2} \sigma^2 S^2 \frac{\partial^2 C}{\partial S^2} dt \end{aligned}$$

Now let  $\Delta = \frac{\partial C}{\partial S}$ , and we get the Black Scholes PDE by cancelling out  $dt$  on both sides:

$$(4.1) \quad r \left( C - S \frac{\partial C}{\partial S} \right) = \frac{\partial C}{\partial t} + \frac{1}{2} \sigma^2 S^2 \frac{\partial^2 C}{\partial S^2}$$

### 5. BLACK SCHOLES EQUATION IN TERMS OF GREEKS

An easy way to intuitively think of the Black Scholes equation is as follows: roughly speaking, the PnL of an option is the sum of  $\theta$  and  $\Gamma$ , note that  $\theta$  is negative and  $\Gamma$  is positive. This PnL should be matched by the growth at the risk free rate of the initial hedged portfolio and hence we get,

$$(5.1) \quad \theta + \frac{1}{2} \sigma^2 S^2 \Gamma = r(C - S\Delta).$$

Writing the same equation with the partial derivatives we get the more well known form of the Black Scholes equation:

$$(5.2) \quad \frac{\partial C}{\partial t} + \frac{1}{2} \sigma^2 S^2 \frac{\partial^2 C}{\partial S^2} = r(C - S \frac{\partial C}{\partial S})$$

Now suppose that you have delta hedged the option. This means the right hand side of Equation 5.1 is 0, i.e.,

$$(5.3) \quad \theta + \frac{1}{2} \sigma_I^2 S^2 \Gamma = 0,$$

where  $\sigma_I$  is the implied volatility of the underlying asset. If we own the option and the realized volatility turns out to be  $\sigma_R$ , then our P&L in a short time period  $\Delta t$

is

$$(5.4) \quad \text{P\&L} = \frac{1}{2} S^2 \Gamma (\sigma_R^2 - \sigma_I^2) .$$

$S^2 \Gamma = S^2 \frac{\partial^2 C}{\partial S^2}$  is known as the *Dollar Gamma* and has the same unit as the price of the option  $C$ , hence the name Dollar Gamma.

The net P&L is given by the equation:

$$(5.5) \quad \text{Net P\&L} = \frac{1}{2} \int_0^T S^2 \Gamma (\sigma_R(t)^2 - \sigma_I^2) dt .$$

## 6. PRICE AND VALUE

According to Derman

Price is simply what you have to pay to acquire a security; value is what it is worth. The price is fair when it is equal to the value.

### MOTIVATION FOR RISK NEUTRAL PRICING

Suppose there is a horse race with two horses. People place bets of total amounts  $x$  and  $y$  on the two horses, respectively. The bookmaker has total deposits amounting to  $x + y$ .

Based on the bets placed, the bookmaker implies and quotes the odds as  $x/y$  for the two horses. (This is artificial, as the odds must be known in advance for the people to bet the money).

The total amount of money the bookmaker has to pay if the first horse wins is  $x + (y/x)x = x + y$ , i.e., return the original amount  $x$  and the odds times the original amount.

Similarly, the total amount of money the bookmaker has to pay if the second horse wins is  $y + (x/y)y = x + y$ .

In both cases the bookmaker has to pay  $x + y$ , which is exactly what he had as a deposit. The bookmaker has no interest in the outcome of the race!

### RISK NEUTRAL PROBABILITY

The Risk neutral probability of a certain event, where the event is described by a financial contract, can be thought of as the *market price probability*, i.e. the probability inferred from the price that the market is willing to pay for that contract.

Let  $B$  be a binary contract, which pays \$1 at time  $T$  if an event  $E$  occurs and nothing if  $E$  does not occur, then the risk neutral probability of  $E$  is:

$$P(E) = \frac{\text{Price(Contract paying \$1 at time } T \text{ if } E \text{ occurs)}}{\text{Price(Contract paying \$1 dollar at time } T \text{ no matter what)}} .$$

**Risk Neutral Price of a Stock.** Let us assume that the price  $S(t)$  of a stock follows geometric Brownian motion. The stochastic differential equation (SDE) followed by the stock price is given by

$$\frac{dS(t)}{S(t)} = \mu dt + \sigma dW_t,$$

where  $W_t \sim N(0, t)$ .

Using Ito's formula we can solve the above SDE and show that for  $u > t$ ,

$$S(u) = S(t) \exp \left( \left( \mu - \frac{1}{2} \sigma^2 \right) (u - t) + \sigma W_{u-t} \right).$$

We can write the log-normal random variable  $S$  as

$$(6.1) \quad S = e^X,$$

where  $X(u)$  is normally distributed:

$$(6.2) \quad X(u) \sim \mathcal{N} \left( \ln S(t) + \left( \mu - \frac{1}{2} \sigma^2 \right) (u - t), \sigma^2 (u - t) \right).$$

In the Black-Scholes world, we obtain the result  $\mu = r$ . However, a lot of people find this result puzzling. How does  $\mu$  get completely eliminated in the final formula for derivative pricing? We try to explain this by applying the basic principle of risk neutral pricing on the most basic if all securities, the stock itself.

Let  $\mathbb{Q}$  be the risk neutral measure, then by the fundamental theorem of asset pricing, the stock price at time  $t$  is the expected value of the stock at any future time  $T > t$ , discounted back to time  $t$ :

$$S(t) = E^{\mathbb{Q}}[S(T)]e^{-r(T-t)}.$$

Using 6.1, we can write the above as

$$\begin{aligned} S(t) &= E^{\mathbb{Q}}[e^{X(T)}]e^{-r(T-t)} \\ S(t) &= e^{\ln S(t) + \left( \mu - \frac{1}{2} \sigma^2 \right) (T-t) + \frac{\sigma^2(T-t)}{2}} e^{-r(T-t)} \\ S(t) &= S(t) e^{(\mu-r)(T-t)} \\ 1 &= e^{(\mu-r)(T-t)}, \end{aligned}$$

and since  $T - t \neq 0$ , the last equation implies,

$$\mu = r.$$

The insight is that the special form of the SDE and its solution coupled with the no arbitrage asset pricing theorem forces the remarkable equality of  $\mu$  and  $r$ .

Note that to help with the computation of  $E^{\mathbb{Q}}[e^X]$ , we can make use of the moment generating function  $\phi_X$  of a normal random variable  $X$  with mean  $m$  and variance  $v^2$ :

$$\phi_X(s) = E[e^{sX}] = e^{ms + \frac{v^2 s^2}{2}}.$$

Evaluating  $\phi_X$  at  $s = 1$ , we get the special result,

$$\phi_X(1) = E[e^X] = e^{m + \frac{v^2}{2}}.$$

In particular, the above identity expresses the mean value of the log-normal random variable  $S = e^X$  in terms of the mean and variance of its underlying normal variable  $X$ .

**Price of a Binary Option.** What is the price of a European binary option which pays \$1 if  $S(T) > K$  and nothing otherwise?

Let  $B(t)$  be the price of the binary option at time  $t < T$ . We have,

$$\begin{aligned}
 B(t) &= E^{\mathbb{Q}}[I(S(T) \geq K)]e^{-r(T-t)} \\
 &= E^{\mathbb{Q}}[I(e^{X(T)} \geq K)]e^{-r(T-t)} \\
 &= E^{\mathbb{Q}}[I(X(T) \geq \ln K)]e^{-r(T-t)} \\
 &= \frac{e^{-r(T-t)}}{\sqrt{2\pi\sigma^2(T-t)}} \int_{\ln K}^{\infty} e^{-\frac{(x-m)^2}{2\sigma^2(T-t)}} dx \\
 &= \frac{e^{-r(T-t)}}{\sqrt{2\pi}} \int_{-\frac{\ln \frac{S(t)}{K} + (r - \frac{1}{2}\sigma^2)(T-t)}{\sigma\sqrt{(T-t)}}}^{\infty} e^{-\frac{z^2}{2}} dz \\
 &= e^{-r(T-t)} \mathcal{N}(d_2),
 \end{aligned}$$

where

$$d_2 = \frac{\ln \frac{S(t)}{K} + (r - \frac{1}{2}\sigma^2)(T-t)}{\sigma\sqrt{(T-t)}}.$$

What is the risk neutral probability that the stock  $S$  at time  $T > t$  is greater than or equal to the strike price  $K$ ?  $\mathcal{N}(d_2)$ .

$$\text{Digital Call: } P^{\mathbb{Q}}(S(T) \geq K) = \mathcal{N}(d_2).$$

Similarly,

$$\text{Digital Put: } P^{\mathbb{Q}}(S(T) \leq K) = 1 - \mathcal{N}(d_2) = \mathcal{N}(-d_2).$$

$$(6.3) \quad \text{Digital Call} + \text{Digital Put} = \text{Riskless Bond}$$

#### PRICE OF A EUROPEAN OPTIONS

The price of a European call option is given by the equation:

$$(6.4) \quad C(t) = S(t)\mathcal{N}(d_1) - e^{-r(T-t)}K\mathcal{N}(d_2),$$

where  $\mathcal{N}$  is the standard normal cumulative distribution function, while

$$(6.5) \quad d_1 = \frac{\ln(S(t)/K) + (r + \sigma^2/2)(T-t)}{\sigma\sqrt{T-t}}, \quad d_2 = \frac{\ln((S(t)/K) + (r - \sigma^2/2)(T-t))}{\sigma\sqrt{T-t}}.$$

Note that as  $t \rightarrow T$ ,  $C_t \rightarrow S_T - K$ .

Also useful is the relationship

$$(6.6) \quad d_2 = d_1 - \sigma\sqrt{T-t}.$$

A more intuitive way to think about  $d_1$  and  $d_2$  is by rewriting them as

$$(6.7) \quad d_1 = \frac{\ln\left(\frac{S(t)e^{r(T-t)}}{K}\right) + \frac{\sigma^2}{2}(T-t)}{\sigma\sqrt{T-t}}, \quad d_2 = \frac{\ln\left(\frac{S(t)e^{r(T-t)}}{K}\right) - \frac{\sigma^2}{2}(T-t)}{\sigma\sqrt{T-t}}.$$

The price of a European put option is given by:

$$(6.8) \quad P(t) = -S(t)\mathcal{N}(-d_1) + e^{-r(T-t)}K\mathcal{N}(-d_2),$$

#### THE PUT-CALL PARITY

For European options, the put-call parity is given by the equation

$$(6.9) \quad S(t) + P(t) - C(t) = Ke^{-r(T-t)}.$$

Another way of remembering the put-call parity is by the phrase: *long call and short put is the same as a forward*.

$$(6.10) \quad C(t) - P(t) = S(t) - Ke^{-r(T-t)}.$$

The put-call parity only holds for European options.

**6.1. Price of a forward contract.** A forward contract is an over the counter (OTC) agreement to buy a stock  $S$  at time  $T$  and price  $K$ . The value  $F$  of this forward contract as a function of time  $t \leq T$  is:

$$(6.11) \quad F(t) = S(t) - e^{-r(T-t)}K.$$

Note that

$$F(t) = 0 \iff K = S(t)e^{r(T-t)}.$$

The economic interpretation of the last equation is simple: it should be free to enter into an at the money forward contract. All other prices will lead to arbitrage.

#### WHY AMERICAN CALLS HAVE THE SAME PRICE AS EUROPEAN CALLS?

An American call option of a stock which pays no dividend has the same price as that of the European option. Let the strike price of the call be  $K$  and its maturity be  $T$ . The optimal strategy for a holder of an American call is to exercise it when the value of the option is the same as its intrinsic value.

At time  $T$ , the payout of the call plus a bond that pays  $K$  at time  $T$  is

$$(S_T - K)^+ + K = \max\{S_T, K\} \geq S_T$$

So at time  $t$ , if we setup a portfolio that consists of the above call and the above bond, then we have to spend

$$X_t = V_t + Ke^{-r(T-t)},$$

where  $V_t$  is the price of the call option at time  $t$ . At time  $T$ , the value of this portfolio will dominate the stock price  $S_T$ . As a result, no arbitrage implies that at time  $t$ ,

$$X_t > S_t$$

Otherwise, we can short<sup>1</sup> one share of stock at time  $t$ , and use the proceeds to setup this portfolio; at time  $T$ , we have zero probability of losing money, and

---

<sup>1</sup>Short the over-priced asset, go long on the under-priced

have a positive probability  $P(S_T < K)$  of making money. This is an arbitrage! Combining the last two equations, we have that

$$V_t > S_t - Ke^{-r(T-t)} > S_t - K.$$

Which shows that the value of the option is always strictly greater than its intrinsic value  $S_t - K$ , therefore the holder should not exercise this option before its maturity  $T$ .

## 7. VTS: VOLATILITY TRADING STRATEGIES

### 7.1. Volatility, Skew, Smile and the Vol Surface.

- The volatility  $\sigma$ , commonly called vol, is usually expressed on an annualized basis. We can therefore assume that the unit of volatility is per square root of time. Similarly, the unit of variance (var) is per unit of time. Correspondingly, the quantities  $\sigma\sqrt{T}$  and  $\sigma^2T$  are dimensionless quantities and represent the total vol and total var attained in time  $T$ .
- Vol Surface: Volatility  $\sigma$  as a function of strike  $K$  and maturity  $T$ ,  $\sigma(K, T)$ .
- ATMF: At the money forward.
- Normalized Strike (NS): For a given maturity  $T$  and strike  $K$ , if the ATMF vol is  $\sigma_{ATMF}$  and the forward value of the underlying asset is  $F$ , then the normalized strike  $NS$  is defined by the equation,

$$NS = \frac{1}{\sigma_{ATMF}\sqrt{T}} \log \left[ \frac{K}{F} \right].$$

- Vol Skew: For a given maturity  $T$ , the slope of the vol surface slice, defined by the equation,

$$Skew = -\frac{1}{\sigma_{ATMF}} \frac{d\sigma(NS)}{dNS}.$$

- Smile: For a given maturity  $T$ , the curvature of the vol slice defined by the equation,

$$Smile = -\frac{1}{\sigma_{ATMF}} \frac{d^2\sigma(NS)}{dNS^2}.$$

- Vol Term Structure: Vol as a function of time to maturity,  $\sigma(T)$ .
- In the Black Scholes world,  $\sigma(K, T)$  is a constant.
- In the FX options market, the vol smile actually looks like a human smile :).
- In the Equity derivatives market, the vol smile is a smirk with a downward slope.
- Sticky by Strike Vol: The vol does not change as a function of the spot  $S$ . The vol still changes as a function of the strike  $K$ . For example, if at the money vol is  $\sigma_0$ , one sticky by strike vol model can be

$$(7.1) \quad [TODO]\sigma(S, K, T) = \sigma(K, T) - b(K - S_0)$$

**7.2. A Short Summary of the Local Vol Model.** The Local Vol model assumes one factor geometric Brownian motion for the underlying asset where the volatility is a deterministic function of spot and time. Crucially, the following are also assumed to be deterministic:

- Interest rates
- Equity funding spreads
- Dividend yield

The aim of the model is to match the non-arbitrageable input implied vol surface at all strikes and maturities. The model can be thought of as a very good interpolator of the implied volatility surface, and allows us to accurately price European styles payoffs.

The local vol surface is analogous to the forward rates. Given two zero coupon bonds, with maturities  $T_1$  and  $T_2$ , such that  $T_1 < T_2$ , the forward rate  $r$  of a zero coupon starting at time  $T_1$  and maturing at time  $T_2$  is given by the equation

$$r_2 T_2 = r_1 T_1 + r(T_2 - T_1).$$

This forward rate  $r$  is the only rate that is consistent with our market rates  $r_1$  and  $r_2$ .

Similarly, given vols of various maturities and strikes (the discrete vol surface formed by market quotes), the *forward vol* aka the local vol is the vol surface which is consistent with the existing discrete vol surface of market quotes.

**7.3. A Short Summary of the Risky Log-OU model.** The Risky Log-OU enhances the local vol model in that it no longer assumes a deterministic funding spread. The model simulates the funding spread of an entity as a stochastic process. The aim of this model is to capture the correlation between equity and funding spreads, giving us conditional (path-wise) risky discounting based on equity levels, making it a richer model compared to a simple local volatility model which has fixed risky discount factors.

#### 7.4. General Rules.

- There are four basic *reserves*
  - (1) Vega
  - (2) Skew
  - (3) Funding Spreads
  - (4) Dividend
- *Funding Spreads*: Let  $r_C(t)$  be the CSA backed short rate, i.e., the agreed overnight rate paid on collateral according to the CSA. We now consider an asset, whose price is some process denoted by  $S(t)$ . Let  $r_R(t)$  be the repo rate of this asset, i.e., the short interest rate we can get if we use the asset as the collateral. Piterbarg calls the difference  $r_R(t) - r_C(t)$  as the stock lending fee [6]. Following Piterbarg [6], we also define the short rate



for unsecured funding by  $r_F(t)$ . Now comes the funding spread, which is defined as [6]

$$s_F(t) = r_F(t) - r_C(t).$$

In essence funding spread of an entity represents the market view of credit default risk of that entity.

In the context of equities, *Stock lending fee* is the extra interest rate charged on top of say LIBOR (or the CSA interest rate) if one wants to borrow money to buy a stock. One typically puts the stock as the collateral for the funding, which is a risky asset and therefore the lender wants to be rewarded with a risk premium for taking this risk. This is reflected via the stock lending fee.

It is expected that

$$(7.2) \quad r_C(t) \leq r_R(t) \leq r_F(t)$$

*Negative Funding Spreads:* Funding spreads can be negative. This can especially happen for short dated maturities. If we have very good credit rating, a short term loan can be obtained at a rate very close to LIBOR. Suppose we now use the cash we obtained to buy a stock and then lend the stock in the repo market. Effectively, this makes our borrowing cost lower than LIBOR, i.e., we have a negative funding spread.

- One wrong argument goes like this. Call prices should go down as funding spreads increase. This is because the large funding spreads are an indication that the equity is being considered very risky. After all, as a lender, if the money you gave away has the stock as the collateral, the spread you charge is proportional to the risk you imagine for the stock. This is also the reason that Put prices go up as funding spreads increase.

The above argument is wrong. In reality, call prices go up as funding spreads go up and put prices go down as funding spreads go up. The reason – option pricing is done via hedging – not via the market view argument presented above. If you are selling a call option, you hedge it with buying delta. The cost of funding to buy your delta increases as the funding spreads increase and as a result, the cost of the option you are selling goes up as well. Similarly, if you are selling a put option, you are going to short some delta. However, when you short a stock, you have to pay a borrowing fee which moves in opposite direction to the funding spread [check this last sentence].

- Call prices go down as dividends increase. The equity spot will be down if dividends increase which results in a lower call price. Correspondingly, put prices go up if dividends go up.

## 7.5. Miscellaneous.

- *Call Overwrite (Buy-Write):* Go long the stock and short a slightly OTM short dated Call. This caps your profit but the proceeds of the premium contribute towards reducing the cost of stock buying and therefore enhance the yield.

### 7.6. Out-performance Options.

- An out-performance option of an asset  $A$  over another asset  $B$  pays  $\max\{A - B, 0\}$  at expiration. By convention, out-performance options are always quoted as a *Call of A over B*.
- Out-performance options are usually European.
- Out-performance options are short correlation. Extreme case: suppose you have two assets,  $X$  and  $Y$ :

$$(7.3) \quad \sigma_{X-Y}^2 = \sigma_X^2 + \sigma_Y^2 - 2\rho_{XY}\sigma_X\sigma_Y,$$

If  $\rho_{XY}$  goes up to 100% the variance on the left is minimum, i.e. the option has the least price. If  $\rho_{XY}$  goes all the way down to  $-100\%$ , the variance  $\sigma_{X-Y}^2$  is maximized.

**7.7. Worst-Of and Best-Of Options.** We will use  $W$  for worst-of and  $B$  for best-of, and the subscripts  $C$  and  $P$  for Calls and Puts respectively. Assume there are  $n$  underlying assets. Then, the payoff worst-of and best-of options is given by

$$(7.4) \quad W_C = \max\left\{\min_{i=1,2,\dots,n}\{S_i - K\}, 0\right\},$$

$$(7.5) \quad W_P = \max\left\{\max_{i=1,2,\dots,n}\{K - S_i\}, 0\right\},$$

$$(7.6) \quad B_C = \max\left\{\max_{i=1,2,\dots,n}\{S_i - K\}, 0\right\},$$

$$(7.7) \quad B_P = \max\left\{\min_{i=1,2,\dots,n}\{K - S_i\}, 0\right\}.$$

Worst-of call options are traded much more than best-of call options because worst-of call options are very cheap.

Worst-of put options are traded much more than best-of put options despite being quite expensive. This is due to the fact that worst-of put options provide protection and are high in demand.

#### 7.7.1. Worst-Of Calls.

- If we lower the correlation between the assets, the price of worst-of call decreases. Intuitive reason, the stocks are more dispersed and the chance that one of them is below the strike is relatively higher. Indeed, worst-of call will be worthless even if a single asset is below the strike. Hence, worst-of call have a lower price if the correlation is low.
- If the correlation goes up, the price of worst-of call also goes up. Extreme case, if we have two assets which are 100% correlated, the price of the worst-of call on the two assets is the same as the price of cheaper of the two.
- Related to the above, a worst-of call is long correlation.
- Typically, clients buy worst-of calls as it gives them a cheaper way of getting exposure to the upside. As a result of the trade, the clients go long correlation. When a trading desk sells worst-of calls, it gets a short correlation exposure, typical position of an exotic derivative desk.

### 7.7.2. *Worst-Of Puts.*

- Worst-of put options are very expensive but still traded as they provide protection. The price of a worst-of put is higher than any of the puts on the underlying assets.
- If we lower the correlation between the assets, the price of worst-of put increases. Intuitive reason, the stocks are more dispersed and the chance that one of them is deep down below the strike is relatively higher. Indeed, worst-of put will pay according to the worst asset which has gone down the most. Hence, worst-of puts have a higher price if the correlation is low.
- If the correlation goes up, the price of worst-of puts goes down (but only relatively, they are already quite expensive). Extreme case, if we have two assets which are 100% correlated, the price of the worst-of puts on the two assets is the same as the price of more expensive of the two.
- Related to the above, worst-of puts are short correlation.
- When we sell worst-of puts, we are long correlation, i.e. we want the correlation to go up so that we have to pay less on the short put positions.
- A trading desk usually buys worst-of-puts (clients like to sell these since they look expensive when compared to basket puts, or individual puts). The desk therefore ends up short correlation (classic exotic position). Due to correlation skew the Bid price can be quite high compared to other buyers in the market. Therefore, there is sometimes pressure to bid lower to be competitive since other market participants might not be charging as much correlation skew.

7.7.3. *Worst-Of Options and the Correlation Skew.* Suppose we are long a WO put, which gives us a short correlation exposure. If the spot goes down, the correlation will go up and because of the correlation skew, the correlation is likely to go up high enough such that the increase in value of our put due to spot going down will be offset to a great extent by the correlation going up.

On the other hand, if the spot goes up the correlation goes down, but because of the correlation skew, it goes down very little. Our put goes cheaper due to the spot going up and the offsetting benefit we get from the correlation going down is small due to the correlation skew.

In either case, the correlation skew is hurting us. Hence, if we are buying a WO put, we should charge extra for the correlation skew.

7.8. **Convexity.**  $9 = 5^2 - 4^2 > 4^2 - 3^2 = 7$  and this is called convexity. If you are long a convex payoff, you make more money on the up move than you would lose on a down move of equal magnitude.

Mathematically, for  $x, \delta x > 0$ ,  $(x + \delta x)^2 - (x - \delta x)^2 \simeq 4x\delta x > 0$ .

## DERIVATIVES ON FOREIGN INDICES

Let  $X$  be the value of 1 JPY in GBP as a function of time. We also suppose that  $X_0$  is the value of  $X$  at  $t = 0$  and  $X_T$  is the value of  $X$  at expiration time  $T$ . Let

$K$  be the strike price in JPY and let  $S$  be the price of a stock in JPY, again as a function of time.

The underlying asset in all of the following cases is a stock with the value  $S$ , denominated in JPY.

- (1) *Foreign-market derivatives*: Buy JPY denominated derivative by converting your GBP into JPY today. At expiry, you will convert the JPY payoff back into GBP at the prevailing exchange rate. You are exposed to both the asset and the FX.

$$(7.8) \quad \max\{0, X_T(S_T - K)\}.$$

- (2) *Compo*: The derivative in this case derives its value by converting the value of the Japanese asset from JPY to GBP at the prevailing exchange rate. You are still exposed to the asset and the FX but this exposure is different from the previous case.

$$(7.9) \quad \max\{0, X_T S_T - K_{\$}\}.$$

- (3) *Quanto*: The payoff of the derivative in this case is the JPY value of the derivative at the time of expiry converted to GBP at a guaranteed exchange rate decided at the start of the contract. You are still exposed to the performance of the underlying asset.

$$(7.10) \quad \max\{0, X_0(S_T - K)\}.$$

## 8. CREDIT

### 8.1. Glossary.

- **Bond Yield**: A single discount number, under which the sum of the present values of all the cash flows of a bond equal its market price. Let  $t$  be the current time,  $C_i$  the cash flow generated by the bond at time  $t_i \geq t$ ,  $P(t)$  the current market price of the bond. Then the bond yield  $y$  is defined by the equation,

$$(8.1) \quad P(t) = \sum_i C_i e^{-y(t_i - t)}, \quad t \leq t_i.$$

- **Par Value**: The bond's principal, also known as the face value.
- **Par Yield**: The coupon rate that causes the bond price to match its Par Value.
- **Strip(s)**: Zero coupon bonds that are synthetically created by selling or buying the coupon of a treasury bond separately from the principal.
- **Spread**: Spread is the constant (absolute) shift to the zero-coupon discount curve in all scenarios that is required to ensure that the model value of the bond (average value over all scenarios) equals the observed market price [3].
- **Rally**: When bond prices go up and yields go down.
- **Sell-Off**: When bond prices go down and yields go up (opposite of a rally).

- Long Credit using CDS: When you have sold protection via a credit default swap. In this case you are long the credit of the company you have sold protection on. You are short a put option on the company's debt. You will make money if the CDS spread tightens.
- Short Credit using CDS: When you have bought protection via a credit default swap. In this case you are short the credit of the company you have bought protection on. You are long a put option on the company's debt. You will make money if the CDS spread widens.
- DTS: Spread duration times spread. This is a measure of credit risk which gives better estimate of credit volatility than using spread duration. If  $s$  is the spread and  $D_s$  the spread duration of a credit security, the spread return  $r_{spread}$  can be approximated as

$$(8.2) \quad r_{spread} \simeq -D_s \Delta s,$$

or equivalently, after multiplying and dividing by  $s$ , as,

$$(8.3) \quad r_{spread} \simeq -D_s \times s \left( \frac{\Delta s}{s} \right).$$

Empirically, the markets have shown the last equation above is a better approximation, especially for volatility prediction of  $r_{spread}$  and hence  $D_s \times s$  (DTS) is a better measure of credit risk.

- Bond Price–Yield–Duration–Convexity: Let  $y$  be the yield of a bond with price  $P$ . Then

$$(8.4) \quad \frac{1}{P} dP \approx -D dy + \frac{1}{2} C dy^2,$$

where the duration  $D$  and the convexity  $C$  are defined as

$$D = -\frac{1}{P} \frac{dP}{dy},$$

$$C = \frac{1}{P} \frac{d^2 P}{dy^2}.$$

- Negative Convexity: For non-callable bonds, the negative relationship between price and yield is usually convex. However, the price-yield relation for callable bonds is typically concave due to the callability feature. This phenomenon is called negative convexity.
- Duration sensitivity: Rate of change of duration with respect to yields (or interest rates).

$$\begin{aligned} \frac{dD}{dy} &= \frac{d}{dy} \left( -\frac{1}{P} \frac{dP}{dy} \right) \\ &= \frac{1}{P^2} \frac{dP}{dy} \frac{dP}{dy} - \frac{1}{P} \frac{d^2 P}{dy^2} \\ &= D^2 - C. \end{aligned}$$

Therefore, the rate of change of duration equals duration squared minus the convexity.

For bonds with positive convexity, the duration can increase or decrease depending on the interplay between the  $D^2$  and the  $C$  term.

For bonds with negative convexity, the duration increases with increase in interest rates and vice-versa.

Historically, the US aggregate index has shown rates-up-duration-up behaviour due to the callability feature of various bonds in the index composition. This pattern holds even more strongly in the US Mortgage Backed Securities (MBS). The US Credit Index on the other hand has historically followed the more intuitive rates-up-duration-down behaviour. Also note that there are indices which do not show a reliable impact of rate changes on duration changes either way.

- **CDS Spread:** Credit default swaps are priced in terms of a spread usually expressed in basis points of the notional value. A CDS quote of 4.55 means that the CDS is pricing at a spread of 4.55 bp, or .0455% i.e., to buy \$10000 of protection, you have to pay \$4.55 per year.
- **EDF:** Expected Default frequency. The default probability predicted by Moody's KMV model.
- **Defensiveness:** In general, defensiveness means the negative correlation of a signal with down markets. In the context of credit, defensiveness of a signal means the negative correlation of the signal with the expected default frequency.
- **The economic cycle and bond prices (The Economist):** Bond prices move in the opposite direction to confidence in the market; bond yields go in the same direction as confidence. When the outlook for the economy is bleak, yields fall sharply as investors rush to the safety of bonds. As the outlook brightens, bond prices start to fall and yields start to rise again. Bond prices are thus countercyclical most of the time. This feature makes them very attractive diversifiers for equities, the prices of which are more procyclical, moving up and down in tandem with the economic cycle.

**8.2. Hazard Rate.** The hazard rate (also called default intensity) is defined as a number  $h$  such that the probability of default in a certain time interval  $[t, t + \Delta t]$ , *conditional* on no earlier default is given by  $h\Delta t$ . In general, the hazard rate will be different when different time intervals are considered, and in those cases  $h(t)$  is defined as a time dependant hazard rate.

Let  $X$  be the random variable representing the time (in years) when the company defaults. For a simple hazard rate model, we assume that  $h$  is the average hazard rate and the time to default is an exponentially distributed random variable  $X$ :

$$P(X \leq t) = 1 - \exp(-ht),$$

Note that  $h$  is playing the same role that the more familiar  $\lambda$  plays in a classically defined exponential random variable  $X$ , with pdf:

$$f_X(x) = \lambda \exp(-\lambda x),$$

and CDF:

$$F_X(x) = 1 - \exp(-\lambda x)$$

The interpretation of  $h$  is given by the equation,

$$\begin{aligned}
 P(t \leq X \leq t + \Delta t | X > t) &= \frac{P(t \leq X \leq t + \Delta t \text{ and } X > t)}{P(X > t)} \\
 &= \frac{P(t \leq X \leq t + \Delta t)}{P(X > t)} \\
 &= \frac{(1 - \exp(-h(t + \Delta t))) - (1 - \exp(-ht))}{\exp(-ht)} \\
 &= 1 - \exp(-h\Delta t) \\
 &\approx h\Delta t.
 \end{aligned}$$

In other words, for this model, the hazard rate (or the default intensity) determines that the probability of default in a small time interval  $\Delta t$  is approximately  $h\Delta t$

Note that

- (1) The default probabilities backed out of bond prices or credit default swap spreads are risk-neutral default probabilities.
- (2) The default probabilities backed out of historical data are real-world default probabilities.

**8.3. Put-Call Parity for the Merton Model.** The fundamental balance sheet equation of a firm is given by

$$V(t) = E(t) + D(t),$$

that is, at any time  $t$ , the firm's assets  $V(t)$  are a sum of the firm's equity  $E(t)$  and the firm's liabilities (or Debt)  $D(t)$ .

The Merton model shows that the firm's equity is a call option on the firm's assets, expiring at time  $T$ , having a strike  $F$ :

$$E(t) := (V(T) - F)^+ = \max\{V(T) - F, 0\}.$$

The strike  $F$  is implied by the initial value of the debt  $D(0)$  and a risky interest rate  $k_D$ :

$$F := D(T) = D(0) \exp(k_D T) = D(t) \exp(k_D (T - t))$$

.

The put-call parity implies:

$$V(t) + P(t) - C(t) = Fe^{-r(T-t)},$$

where  $r$  is the risk-free rate. Since  $C(t) = E(t)$ , using the balance sheet equation gives  $V(t) - C(t) = D(t)$  and substituting this in the put-call parity above gives us:

$$D(t) + P(t) = Fe^{-r(T-t)}.$$

Now use the value  $F = D(t)e^{k_D(T-t)}$ :

$$\begin{aligned} D(t) + P(t) &= D(t)e^{k_D(T-t)}e^{-r(T-t)} \\ P(t) &= D(t) \left( e^{(k_D-r)(T-t)} - 1 \right) \\ \frac{1}{T-t} \ln \left( 1 + \frac{P(t)}{D(t)} \right) &= k_D - r \end{aligned}$$

which implies that the *credit spread* is

$$k_D - r = \frac{1}{T-t} \ln \left( 1 + \frac{P(t)}{D(t)} \right).$$

## 9. PROBABILITY AND STOCHASTIC CALCULUS

**Definition 1** (Probability Space). . A triple  $(\Omega, \mathcal{F}, \mathcal{P})$ .  $\Omega$  is a set,  $\mathcal{F}$  is a sigma-algebra on  $\Omega$  and  $\mathcal{P} : \mathcal{F} \rightarrow [0, 1]$ , such that

- (1)  $\mathcal{P}(\Omega) = 1$
- (2) For any countable union of mutually disjoint sets in  $\mathcal{F}$ , the function  $\mathcal{P}$  is additive.

**Definition 2** (Probability Measure). The function  $\mathcal{P}$  above is called a probability measure.

**Definition 3** (Random Variable). A random variable is a  $\mathcal{P}$ -measurable function  $X : \Omega \rightarrow \mathbb{R}$ .

**Conditional Probability Notation.** Let  $X$  and  $Y$  be random variables with a joint distribution. The expression  $Y|X$  is not a random variable. It is simply a notation which dictates that any operation on the random variable  $Y$  must be done so using the conditional distribution of  $Y$  given  $X$ . That is,  $X$  should be treated as a known constant. Using this notation, we write, for example, the conditional expectation of  $Y$  given  $X$  as  $E(Y|X)$ . This conditional expectation is indeed a random variable as it is a function of  $X$ . To evaluate this conditional expectation, we use the conditional density of  $Y$  given  $X$  and sum or integrate over all values of  $Y$  [1]:

$$E(Y|X) = \int y f_{Y|X}(y) dy.$$

**9.1. Modes of convergence of a sequence of random variable.** Let  $X_n$  be a sequence of random variables defined on a probability space.

We say that

- (1)  $X_n \rightarrow X$  almost surely if for sufficiently large  $n$   $P(|X_n(\omega) - X(\omega)| < \epsilon) = 1$ .
- (2) Mean square convergence, which is stronger than convergence in probability. [TODO]
- (3)  $X_n \rightarrow X$  in probability if
- (4)  $X_n \rightarrow X$  in distribution if for all  $z \in \mathbb{R}$ ,  $F_n(z) \rightarrow F(z)$ , where  $F_n$  is the distribution function of  $X_n$  and  $F$  is distribution function of  $X$ .



For a Venn diagram of modes of convergence, see Figure 7-5 of [5].

**9.2. Transformation of Random Variables.** Let  $X$  and  $Y$  be two random variables with joint density function  $f_{XY}$ . Given,

$$\begin{aligned} U &= g(X, Y), \\ V &= h(X, Y), \end{aligned}$$

what is the joint density of  $f_{UV}$ ?

We assume that we can *invert* the transformation and express  $X$  and  $Y$  as,

$$\begin{aligned} X &= \phi(U, V), \\ Y &= \psi(U, V). \end{aligned}$$

Secondly,

$$dA = dxdy = |J(x, y)| dudv,$$

where,

$$J(x, y) = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix}.$$

Then,

$$\begin{aligned} P((U, V) \in A) &= \int_A f_{UV}(u, v) du dv \\ &= P((X, Y) \in B) \\ &= \int_B f_{XY}(x, y) dx dy \\ &= \int_B f_{XY}(\phi(u, v), \psi(u, v)) dx dy \\ &= \int_B f_{XY}(\phi(u, v), \psi(u, v)) |J(x, y)| du dv. \end{aligned}$$

Therefore,

$$f_{UV}(u, v) = f_{XY}(\phi(u, v), \psi(u, v)) |J(x, y)|.$$

**Brownian Motion.** Brownian motion is a stochastic process  $B(t)$ , such that

- (1)  $B(0) = 0$ .
- (2) For any  $t_1 < t_2$ ,  $B(t_2) - B(t_1)$  is independent of  $B(t_1)$  and normal with zero mean and variance  $t_2 - t_1$ .
- (3)  $B(t)$  is continuous almost surely.

**Stopping Time.** Intuitively speaking, a stopping time is the time at which a stochastic process satisfies a certain rule. However, for a stopping time, the rule must be defined in a way that at any given time it may be tested by looking at only the past and present values of the stochastic process. For example: buy Microsoft as soon as the stock price goes below \$100 is a valid rule for a stopping time. On the other hand: sell Microsoft first thing in the morning if the closing price that day is less than \$50 per share is not a stopping time, since in the morning we can not tell what closing price will be attained at the end of the day.

Let  $\tau_m := \inf\{t : B(t) = m\}$ , i.e., the random variable  $\tau_m$  is the time when Brownian motion  $B$  hits the level  $m$  for the first time.

Notice that

$$(9.1) \quad \tau_m = \inf\{t : B(t) \geq m\}$$

What is the probability  $P(\tau_m \leq T)$ ?

$$P(\tau_m \leq T) = P(\tau_m \leq T \cap B(T) \geq m) + P(\tau_m \leq T \cap B(T) < m).$$

Now,

$$(9.2) \quad P(\tau_m \leq T \cap B(T) \geq m) = P(B(T) \geq m),$$

and by using the reflection principle,

$$(9.3) \quad P(\tau_m \leq T \cap B(T) < m) = P(B(T) \geq m),$$

Therefore,

$$(9.4) \quad P(\tau_m \leq T) = 2P(B(T) \geq m).$$

$$(9.5) \quad P(\tau_m \leq t) = 2P(B(t) \geq m)$$

$$(9.6) \quad = 2 \frac{1}{\sqrt{2\pi}} \int_m^\infty e^{-\frac{x^2}{2t}} dx$$

$$(9.7)$$

Therefore, the density function of  $\tau_m$  is obtained by differentiating the last expression with respect to  $t$ :

$$(9.8) \quad f_{\tau_m}(t) =$$

**9.3. Copula.** A copula is a distribution function  $C : \mathbb{R}^n \rightarrow [0, 1]$  such that all the marginals are standard uniform random variables.

By definition,

$$C(1, 1, \dots, 1, u_j, 1, \dots, 1) = u_j.$$

The probability density function  $c$  associated with  $C$  is given by the usual formula,

$$c = \frac{\partial^n C}{\partial u_n \partial u_{n-1} \dots \partial u_1}$$

**9.4. Probability Integral Transform.** Let  $X \sim F$ . Then  $Y = F(X)$  is uniform.

*Proof.*

$$\begin{aligned} P(Y \leq y) &= P(F(X) \leq y) \\ &= P(X \leq F^{-1}(y)) \\ &= F(F^{-1}(y)) \\ &= y. \end{aligned}$$

□

**9.5. Probability Quantile Transform.** Let  $U$  be uniform. Then,  $X = F^{-1}(U) \sim F$

*Proof.*

$$\begin{aligned} P(X \leq x) &= P(F^{-1}(U) \leq x) \\ &= P(U \leq F(x)) \\ &= F(x). \end{aligned}$$

□

**Log-normal random variable.** If the log of a random variable  $X$  is Normal, then  $X$  is called a log-normal random variable.

$$(9.9) \quad \ln X = \sigma Z + \mu,$$

where  $Z$  is a standard normal:  $Z \sim \mathcal{N}(0, 1)$ .

$$(9.10) \quad X = e^{\sigma Z + \mu}.$$

The expected value of  $X$  is given by the formula

$$(9.11) \quad E[X] = e^{\sigma^2/2 + \mu}.$$

**9.6. Ito's Formula.**

**9.7. Ito's Formula for Brownian Motion.** Let  $W(t)$  be a Brownian motion. Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a function with continuous second derivative, then

$$(9.12) \quad df(W(t)) = f'(W(t))dW(t) + \frac{1}{2}f''(W(t)), \quad (\text{Differential Form})$$

or

$$(9.13) \quad f(W(t)) - f(W(0)) = \int_0^t f'(W(s))dW(s) + \frac{1}{2} \int_0^t f''(W(s))dW(s), \quad (\text{Integral Form}).$$

**9.8. Ito's Formula for a useful stochastic process.** Let  $W(t)$  be a Brownian motion. Let us define the stochastic process  $X(t)$  via the stochastic differential equation

$$(9.14) \quad dX(t) = \mu(X(t), t)dt + \sigma(X(t), t)dW(t).$$

We also assume that  $f(t, x)$  is a function with continuous derivatives  $f_t, f_x$  and  $f_{xx}$ . Then, the differential form of Ito's formula expresses the differential increment of  $f(X(t), t)$  by

$$(9.15) \quad df(t, X(t)) = f_t(t, X(t))dt + f_x(t, X(t))dX(t) + \frac{1}{2}f_{xx}(t, X(t))dX(t)^2,$$

while the integral form is

$$(9.16) \quad f(T, X(T)) - f(0, X(0)) = \int_0^T f_t(t, X(t))dX(t) + \int_0^T f_x(t, X(t))dX(t) + \frac{1}{2} \int_0^T f_{xx}(t, X(t))dX(t).$$

For Geometric Brownian Motion (GBM), the SDE can be obtained by setting  $\mu(X(t), t) = \mu_0 X(t)$  and  $\sigma(X(t), t) = \sigma_0 X(t)$  in Equation (9.14). This is usually written as (abusing notation a little bit by identifying  $\mu$  and  $\sigma$  as constants,

$$(9.17) \quad dX(t) = \mu X(t)dt + \sigma X(t)dW(t),$$

where  $X(t)$  is the price of the stock.

The solution of the above SDE is

$$(9.18) \quad X(t) = X(0) \exp \left[ \left( \mu - \frac{1}{2}\sigma^2 \right)t + \sigma W(t) \right]$$

$W(t)$  is a martingale. Indeed, for any  $t > s$ ,

$$\begin{aligned} E[W(t+s)|\mathcal{F}_s] &= E[W(t+s) - W(s) + W(s)|\mathcal{F}_s] \\ &= W(s) + E[W(t+s) - W(s)|\mathcal{F}_s] \\ &= W(s). \end{aligned}$$

$e^{\sigma W(t)}$  is not a martingale. One can check this by applying the Ito's Lemma. Let  $f(W(t)) = e^{\sigma W(t)}$ , then

$$(9.19) \quad df(W(t)) = \sigma W(t)dW(t) + \frac{1}{2}\sigma^2 dt.$$

Because of the non-zero drift term,  $f$  is not a martingale.

$e^{(\sigma W(t) - \frac{1}{2}\sigma^2 t)}$  is a martingale. Again, we can check this using Ito's Lemma. Let  $f(W(t)) = e^{\sigma W(t) - \frac{1}{2}\sigma^2 t}$ , then

$$(9.20) \quad df(W(t)) = -\frac{1}{2}\sigma^2 W(t)dt + \sigma W(t)dW(t) + \frac{1}{2}\sigma^2 W(t)dt = \sigma W(t)dW(t).$$

There is no drift term, hence  $f$  is a martingale.

The moment generating function of  $X \sim N(\mu, \sigma)$  is

$$(9.21) \quad \phi(s) = E[e^{sX}] = e^{s\mu} e^{\frac{s^2\sigma^2}{2}}$$

Using the above result, we see that

$$\begin{aligned}
 E[e^{\sigma W(t+s)} | \mathcal{F}_s] &= E[e^{\sigma(W(t+s)-W(s)+W(s))} | \mathcal{F}_s] \\
 &= e^{\sigma W(s)} E[e^{\sigma(W(t+s)-W(s))} | \mathcal{F}_s] \\
 &= e^{\sigma W(s)} e^{\frac{1}{2}\sigma^2 t^2} \\
 &= e^{\sigma W(s) + \frac{1}{2}\sigma^2 t^2} \\
 &\neq e^{\sigma W(s)}.
 \end{aligned}$$

Therefore,  $e^{\sigma W(s)}$  is not a martingale. Note that we have used the fact that

$$(9.22) \quad W(t+s) - W(s) \sim N(0, t),$$

and also the formula for the moment generating function for  $N(0, t)$ .

**Change of measure.** Consider  $X$ , a standard normal random variable defined on the probability space  $(\mathbb{R}, \mathcal{F}, \mu_X)$ , where  $\mathcal{F}$  is the  $\sigma$ -algebra of open sets in  $\mathbb{R}$  and for a Borel set  $B$ ,

$$\mu_X(B) := \frac{1}{\sqrt{2\pi}} \int_B e^{-\frac{x^2}{2}} dx.$$

The expectation of a Borel measurable function  $f$  under the probability measure  $\mu_X$  is given as

$$E_0[f(X)] = \int_{-\infty}^{\infty} f(X) d\mu_X,$$

where the subscript 0 in the expectation is to emphasize that the mean of  $X$  is 0. What if we now want to change the mean of the distribution to a new number say  $m$ . How is the expectation of  $f$  under the old measure related to the expectation under the new measure? The answer is easy in this case:

$$\begin{aligned}
 E_0[f(X)] &= \int_{-\infty}^{\infty} f(X) d\mu_X \\
 &= \int_{-\infty}^{\infty} f(X) e^{-X^2/2} dX \\
 &= \int_{-\infty}^{\infty} f(X) e^{-mX + \frac{m^2}{2}} e^{-\frac{(X-m)^2}{2}} dX \\
 &\equiv E_m \left[ f(X) e^{-mX + \frac{m^2}{2}} \right].
 \end{aligned}$$

The collection of theorems that tell us how to make drift disappear is commonly called Girsanov theory [7, Ch. 13].

## 10. TIME SERIES ANALYSIS

White noise is our basic building block for time series. We assume that  $\epsilon_t$  is white noise:

$$\epsilon_t \sim \text{i.i.d } N(0, \sigma^2).$$

This assumption implies:

$$\begin{aligned} E(\epsilon_t) &= E(\epsilon_t | \epsilon_{t-1}, \epsilon_{t-2}, \dots) = 0, \\ \gamma_h(t+h, t) &= \text{Cov}(\epsilon_{t+h}, \epsilon_t) = E(\epsilon_{t+h}\epsilon_t) = 0, \quad h \neq 0 \\ \gamma_0(t, t) &= \text{Var}(\epsilon_t) = \text{Var}(\epsilon_t | \epsilon_{t-1}, \epsilon_{t-2}, \dots) = \sigma^2. \end{aligned}$$

**Stationary Process.** There are various kinds of stationary time series. The most common and of particular interest to us is the *weakly stationary* time series. A weakly stationary time series satisfies the following two conditions:

- (1) It has a stationary mean:  $E(x_t) = \mu$ .
- (2) The covariance is stationary:  $\text{Cov}(x_t, x_s)$  is only a function of  $|t - s|$ .

A random walk is not stationary:

$$x_t = x_{t-1} + \epsilon_t$$

Given a starting value of say  $x_0$ , we have,  $E(x_t) = x_0$ , which is indeed a constant, however,  $\text{Cov}(x_t, x_s) = \min(t, s)\sigma_\epsilon^2$ , which is not a function of  $|t - s|$  only.

*The Auto Regressive Moving Average (ARMA) Model.* Let  $L$  be the lag operator:

$$Lx_t = x_{t-1}$$

Various ARMA processes can be described using the lag operator:

$$\left| \begin{array}{l} AR(1) \\ MA(1) \\ AR(p) \\ MA(q) \\ ARMA(p, q) \end{array} \right| \left| \begin{array}{l} x_t = \phi x_{t-1} + \epsilon_t \\ x_t = \epsilon_t + \theta \epsilon_{t-1} \\ x_t = \sum_{k=1}^p \phi_k x_{t-k} + \epsilon_t \\ x_t = \epsilon_t + \sum_{k=1}^q \theta_k \epsilon_{t-k} \\ x_t = \sum_{k=1}^p \phi_k x_{t-k} \epsilon_t + \epsilon_t + \sum_{k=1}^q \theta_k \epsilon_{t-k} \end{array} \right| \left| \begin{array}{l} (1 - \phi L)x_t = \epsilon_t \\ x_t = (1 + \theta L)\epsilon_t \\ a_p(L)x_t = \epsilon_t \\ x_t = b_q(L)\epsilon_t \\ a_p(L)x_t = b_q(L)\epsilon_t \end{array} \right|$$

If  $a_p(L)$  is invertible, we can convert an  $AR(p)$  process to an  $MA(\infty)$  process. Similarly, if  $b_q(L)$  is invertible, we can convert an  $MA(q)$  process to an  $AR(\infty)$  process.

### 10.1. Impulse Response.

10.2. **DTFT.** The discrete time Fourier transform of a doubly infinite, absolutely summable discrete sequence  $x_n$  is given by,

$$X(e^{i\omega}) = \sum_{n=-\infty}^{\infty} x_n e^{-i\omega n}.$$

$X(e^{i\omega})$  is  $2\pi$ -periodic and for real  $x_n$ ,

$$X(e^{i\omega}) = X^*(e^{-i\omega}).$$

$$x_n = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{i\omega}) e^{i\omega n} d\omega.$$

## 11. PORTFOLIO THEORY

Suppose we have a universe of  $n$  assets, and the returns of these assets are given by a  $n$ -vector random variable  $r = [r_1, r_2, \dots, r_n]^T$ .

We assume that the mean of  $r$  is given by

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix},$$

while the covariance matrix of  $r$  is given by

$$\Sigma = \mathbb{E}[(r - \mu)(r - \mu)^T] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{12} & \sigma_2^2 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \cdots & \cdots & \sigma_{ij} & \cdots \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{1n} & \sigma_{2n} & \cdots & \sigma_n^2 \end{bmatrix}.$$

Note that  $\sigma_{ij} = \rho_{ij}\sigma_i\sigma_j$ .

**A Fully Invested Minimum Variance Portfolio.** The return of our portfolio  $r_p$  is a linear combination of weighted asset returns:

$$r_p = w \cdot r := w^T r,$$

where  $w = [w_1, w_2, \dots, w_n]^T$  is a weight vector such that

$$w \cdot \mathbf{1} = \sum_{i=1}^n w_i = 1.$$

However, we want to choose the weight vector  $w$  which minimizes the variance of our portfolio.

Recall that in general for a matrix  $A$  and a random vector  $X$ ,  $\text{Var}(AX) = A\Sigma A^T$ . Therefore,

$$\text{Var}(r_p) = \text{Var}(w^T r) = w^T \Sigma w.$$

We have the following standard minimization problem:

$$\begin{aligned} &\text{minimize: } \frac{1}{2} w^T \Sigma w, \\ &\text{subject to: } w \cdot \mathbf{1} = 1. \end{aligned}$$

The associated Lagrangian is given by,

$$L(w, \lambda) = \frac{1}{2} w^T \Sigma w - \lambda(w \cdot \mathbf{1} - 1).$$

Differentiating with respect to  $w$  and setting the result to zero gives,

$$\Sigma w - \lambda \mathbf{1} = 0,$$

which implies that the best weights vector  $w^*$  is given by,

$$w^* = \lambda \Sigma^{-1} \mathbf{1}.$$

Since  $w^*$  should satisfy the constraint as well,

$$\begin{aligned} w^* \cdot \mathbf{1} &= 1, \\ \lambda \mathbf{1}^T \Sigma^{-1} \mathbf{1} &= 1, \\ \lambda &= \frac{1}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}}, \\ w^* &= \frac{\Sigma^{-1} \mathbf{1}}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}}. \end{aligned}$$

The minimum variance is given by,

$$\begin{aligned} \sigma_{min}^2 &= w^{*T} \Sigma w^*, \\ &= \left( \frac{\Sigma^{-1} \mathbf{1}}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}} \right)^T \Sigma \left( \frac{\Sigma^{-1} \mathbf{1}}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}} \right), \\ &= \left( \frac{\Sigma^{-1} \mathbf{1}}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}} \right)^T \frac{\mathbf{1}}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}}, \\ &= \left( \frac{\Sigma^{-1} \mathbf{1}}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}} \right)^T \frac{(\mathbf{1}^T)^T}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}}, \\ &= \left( \frac{\mathbf{1}^T \Sigma^{-1} \mathbf{1}}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}} \right)^T \frac{1}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}}, \\ &= \frac{1}{\mathbf{1}^T \Sigma^{-1} \mathbf{1}} = \lambda. \end{aligned}$$

**Characteristic Portfolio.** As before, we assume a portfolio of  $n$  assets defined by a weight vector  $w \in \mathbb{R}^n$ .

**Definition 4** (Attribute). *Given an asset, an attribute is a mapping from the asset to a real number.*

*For  $n$  assets, an attribute vector  $a \in \mathbb{R}^n$  is formed by defining the  $i^{th}$  component of  $a$  as the attribute of the  $i^{th}$  asset.*

Some examples of an attribute are: asset return, asset beta, asset size, asset market cap, membership in an industrial sector indicated by a binary variable etc.

**Definition 5** (Unit Exposure). *A portfolio is said to have unit exposure to an attribute vector  $a$  if the weights of the portfolio  $w$  are chosen such that  $w \cdot a = 1$ .*

As an example, the minimum variance problem we solved in the previous section had the constraint  $w \cdot \mathbf{1} = 1$ , i.e. the resulting portfolio had unit exposure to the membership attribute. In other words, a fully invested portfolio has unit exposure to the membership attribute.

**Definition 6** (Characteristic Portfolio). *Given an attribute vector  $a$ , the minimum variance portfolio with unit exposure to  $a$  is called the characteristic portfolio of  $a$ .*



$$\begin{aligned} \min_{w \in \mathbb{R}^n} w^T \Sigma w, \\ w \cdot a = 1. \end{aligned}$$

$$\begin{aligned} \Sigma w - l a &= 0, \\ w &= l \Sigma^{-1} a. \end{aligned}$$

Using the constraint,  $w \cdot a = 1$ , we get

$$w = \frac{\Sigma^{-1} a}{a^T \Sigma^{-1} a}$$

### A Minimum Variance Portfolio with Multiple Constraints.

**Theorem 1.** *Let  $w, f \in \mathbb{R}^n$ ,  $\Sigma \in \mathbb{R}^{n \times n}$ ,  $a_i \in \mathbb{R}^n$  and  $c_i \in \mathbb{R}$ ,  $\lambda > 0$ . The solution of the problem:*

$$\operatorname{argmax}_{w \in \mathbb{R}^n} w^T f - \frac{\lambda}{2} w^T \Sigma w,$$

*with  $m$  exposure constraints,*

$$\begin{aligned} w \cdot a_1 &= c_1, \\ w \cdot a_2 &= c_2, \\ &\vdots \\ w \cdot a_m &= c_m, \end{aligned}$$

*$m < n$ , has the form*

$$w^* = b_f w_f + b_1 w_1 + b_2 w_2 + \dots + b_m w_m,$$

*where  $b_f, b_1, b_2, \dots, b_m$  are scalars and  $w_i$  are the weights of the characteristic portfolio of the attribute vector  $a_i$ .*

*Furthermore, assuming that  $a_i$  are linearly independent, we can find  $w^*$  by solving a linear system of equations completely determined by the characteristic portfolio weights  $w_1, w_2, \dots, w_m$ .*

*Proof.* The Lagrangian associated with the problem is given by,

$$L(w, l) = \left( w^T f - \frac{\lambda}{2} w^T \Sigma w \right) - l_1 w_1 \cdot a_1 - l_2 w_2 \cdot a_2 - \dots - l_m w_m \cdot a_m.$$

Differentiating with respect to  $w$  and setting the result equal to zero gives us the optimal weights vector  $w^*$ ,

$$w^* = \frac{1}{\lambda} \Sigma^{-1} (f - l_1 a_1 - l_2 a_2 - \dots - l_m a_m),$$

which implies that we can write  $w^*$  as,

$$w^* = \frac{1}{\lambda} \Sigma^{-1} f - \frac{l_1}{\lambda} \Sigma^{-1} a_1 - \frac{l_2}{\lambda} \Sigma^{-1} a_2 - \dots - \frac{l_m}{\lambda} \Sigma^{-1} a_m.$$

Since the characteristic portfolio of an attribute vector  $a_i$  has the form  $d_i \Sigma^{-1} a_i$  where  $d_i$  is a real number, the last equation above has the form,

$$w^* = b_f w_f + b_1 w_1 + b_2 w_2 + \dots + b_m w_m,$$

for some real numbers  $b_f, b_1, b_2, \dots, b_m$ . This proves the first part of the theorem.

Observe that we can write the exposure constraints as

$$A^T w = A^T W b = c,$$

where,

$$A = \begin{bmatrix} a_1 & a_2 & \dots & a_m \end{bmatrix}, \quad W = \begin{bmatrix} w_1 & w_2 & \dots & w_m \end{bmatrix},$$

and

$$c = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix}, \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}.$$

Note that

$$w = Wb,$$

and the matrix  $W$  can be re-written as,

$$\begin{aligned} W &= \begin{bmatrix} w_1 & w_2 & \dots & w_n \end{bmatrix} \\ &= \begin{bmatrix} d_1 \Sigma^{-1} a_1 & d_2 \Sigma^{-1} a_2 & \dots & d_n \Sigma^{-1} a_n \end{bmatrix} \\ &= \Sigma^{-1} \begin{bmatrix} a_1 & a_2 & \dots & a_n \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \dots \\ d_n \end{bmatrix} \\ &= \Sigma^{-1} A D. \end{aligned}$$

The linear constraints can be written as

$$A^T W b = c$$

$$A^T \Sigma^{-1} A D b = c.$$

since  $A, D, \Sigma^{-1}$  are invertible, a unique solution  $b$  exists.  $\square$

**What is Alpha?** Alpha is a forecast of residual return. Under the Grinold Kahn Framework (GKF), a rule of thumb for alpha is given by:

Alpha = Volatility . IC . score

Mathematically, we can write

$$(11.1) \quad \alpha = \sigma \rho_{z, r_r} z,$$

where  $\rho_{z, r_r}$  (also known as the Information coefficient or IC) is the correlation of score with residual return  $r_r$ ,  $\sigma$  is the volatility of the residual return and  $z$  is the standardized  $z$ -score or the signal value.

**Alpha Derivation.** Let's assume that  $f$  is a forecast vector in  $\mathbb{R}^n$  and  $\Sigma \in \mathbb{R}^{n \times n}$  is a symmetric positive definite covariance matrix and  $\lambda > 0$  is the risk aversion parameter.

Our aim is to solve the following constrained mean-variance optimization problem:

$$\operatorname{argmax}_{w \in \mathbb{R}^n} w^T f - \frac{1}{2} \lambda w^T \Sigma w,$$

subject to constraints of being dollar neutral and neutral to  $k$  other factors  $f_1, f_2, \dots, f_k$ . We can impose these constraints as,

$$[\mathbf{1} \ F]^T w = 0,$$

where  $F = [f_1 \ f_2 \ \dots \ f_k]$  is the  $n \times k$  factor matrix.

Using a factor model based on  $F$ , we run the  $k$ -variate cross-sectional regression on the return vector  $r$ :

$$r = a\mathbf{1} + F\beta_f + \epsilon,$$

where  $\beta_f$  is the unknown  $k$ -vector of regression coefficients.

The above factor model implies that  $\Sigma$ , the covariance matrix of  $r$ , can be decomposed as

$$\Sigma = FV_FF^T + \Delta,$$

where  $V_F$  is the  $k \times k$  covariance matrix of the  $k$  factors and  $\Delta = \operatorname{Var}(\epsilon)$  is a  $n \times n$  diagonal matrix (of specific risk) with positive entries on the diagonal.

This decomposition of  $\Sigma$  together with our factor neutrality constraint  $w^T F = 0$ , implies that  $w^T \Sigma w = w^T \Delta w$ . Therefore, the original optimization problem can be re-written as:

$$\operatorname{argmax}_{w \in \mathbb{R}^n} w^T f - \frac{1}{2} \lambda w^T \Delta w,$$

with the same original constraints

$$[\mathbf{1} \ F]^T w = 0.$$

Forming the Lagrangian, we find that the optimal holdings  $w_*$  are given by

$$w_* = \frac{1}{\lambda} \Delta^{-1} (f - F\ell - l\mathbf{1}),$$

where  $\ell$  is the  $k$ -vector of Lagrange multipliers associated with the factor neutral constraints and  $l$  is the scalar Lagrange multiplier for the dollar neutral constraint.

The single period optimal return  $r_\pi$  of our portfolio corresponding to the optimal holdings  $w_*$  is given by,

$$\begin{aligned}
r_\pi &= w_*^T r \\
&= w_*^T (a\mathbf{1} + F\beta_f + \epsilon) \\
&= w_*^T \epsilon \\
&= \frac{1}{\lambda} (\Delta^{-1} (f - F\ell - l\mathbf{1}))^T \epsilon \\
&= \frac{1}{\lambda} \left( \underbrace{\Delta^{-\frac{1}{2}} (f - F\ell - l\mathbf{1})}_{\text{Risk adj. alpha} := \hat{\alpha}} \right)^T \left( \underbrace{\Delta^{-\frac{1}{2}} \epsilon}_{\text{Risk adj. spec. ret.} := \hat{\epsilon}} \right). \\
r_\pi &= \frac{1}{\lambda} \hat{\alpha}^T \hat{\epsilon}
\end{aligned}$$

We have written the optimal portfolio return  $r_\pi$  as the inner product of the risk adjusted factor neutral forecast  $\hat{\alpha}$  with the risk adjusted factor neutral specific returns  $\hat{\epsilon}$ , scaled by the risk aversion parameter  $\lambda$ .

Since we cross-sectionally z-score  $\hat{\alpha}$ , the in-sample mean  $\bar{\hat{\alpha}} = 0$  and therefore, the sample estimate of  $\text{Cov}(\hat{\alpha}, \hat{\epsilon})$  is simply given by  $\frac{1}{n-1} \hat{\alpha}^T \hat{\epsilon}$ . This allows the approximation,

$$\begin{aligned}
r_\pi &= \frac{1}{\lambda} \hat{\alpha}^T \hat{\epsilon} \\
&\approx \frac{n-1}{\lambda} \text{Cov}(\hat{\alpha}, \hat{\epsilon}) \\
&= \frac{n-1}{\lambda} \text{Corr}(\hat{\alpha}, \hat{\epsilon}) \sigma_{\hat{\alpha}} \sigma_{\hat{\epsilon}}.
\end{aligned}$$

For the variance of  $r_\pi$ , we have

$$\begin{aligned}
\sigma_{r_\pi}^2 &= \frac{1}{\lambda^2} \hat{\alpha}^T \text{Var}(\hat{\epsilon}) \hat{\alpha} \\
&= \frac{1}{\lambda^2} \hat{\alpha}^T \hat{\alpha} \\
&= \frac{1}{\lambda^2} \|\hat{\alpha}\|_2^2,
\end{aligned}$$

since  $\text{Var}(\hat{\epsilon}) = \text{Var}(\Delta^{-\frac{1}{2}} \epsilon) = \Delta^{-\frac{1}{2}} \text{Var}(\epsilon) \Delta^{-\frac{1}{2}} = \Delta^{-\frac{1}{2}} \Delta \Delta^{-\frac{1}{2}} = I$ , the identity matrix.

Since the sample estimate of  $\sigma_{\hat{\alpha}}^2$  is given by  $\frac{1}{n-1} \|\hat{\alpha}\|_2^2$ , we get

$$\begin{aligned}
\sigma_{r_\pi}^2 &= \frac{1}{\lambda^2} \|\hat{\alpha}\|_2^2 \\
\sigma_{r_\pi}^2 &\approx \frac{n-1}{\lambda^2} \sigma_{\hat{\alpha}}^2,
\end{aligned}$$

which allows us to write,

$$\lambda \approx \sqrt{n-1} \frac{\sigma_{\hat{\alpha}}}{\sigma_{r_\pi}}.$$

Substituting this value of  $\lambda$  in the earlier expression for  $r_\pi$  gives us

$$r_\pi = \sqrt{n-1} \text{Corr}(\hat{\alpha}, \hat{\epsilon}) \sigma_{r_\pi} \sigma_{\hat{\epsilon}}.$$

For a tangible financial interpretation, each term of our return may be labelled as:

$$r_\pi = \underbrace{\sqrt{n-1}}_{\text{Breadth}} \underbrace{\text{Corr}(\hat{\alpha}, \hat{\epsilon})}_{\text{Skill}} \underbrace{\sigma_{r_\pi}}_{\text{Traget risk}} \underbrace{\sigma_{\hat{\epsilon}}}_{\text{Opportunity}}.$$

This also allows us to write the information ratio (IR) as

$$IR = \frac{r_\pi}{\sigma_{r_\pi}} = \sqrt{n-1} \text{Corr}(\hat{\alpha}, \hat{\epsilon}) \sigma_{\hat{\epsilon}}.$$

**What is Beta?** Beta of an asset is the regression coefficient obtained when excess returns of an asset are linearly regressed on the excess returns of the market

Let  $r_A$  and  $r_M$  be the *excess returns* of the asset and the market respectively, then

$$(11.2) \quad \beta = \frac{\text{Cov}(r_A, r_M)}{\text{Var}(r_M)}$$

#### LINEAR REGRESSION

What is the best approximation of a random variable  $Y$  by a constant  $r$ ? The answer depends on what we mean by “best”. The most common way is to identify the best approximation as the one with minimum mean squared error (MSE):

$$r_* = \min_r \mathbb{E}[(Y - r)^2],$$

assuming of course that the mean and variance of  $Y$  is finite.

Set the derivative (with respect to  $r$ ) of the MSE to 0 and simplify to get,

$$\begin{aligned} \frac{d}{dr} \mathbb{E}[(Y - r)^2] &= 0, \\ \frac{d}{dr} \int (y - r)^2 f_Y(y) dy &= 0, \\ \int 2(y - r) f_Y(y) dy &= 0, \\ r \int f_Y(y) dy &= \int y f_Y(y) dy, \\ r &= \mathbb{E}[Y]. \end{aligned}$$

Therefore,  $r_* = \mathbb{E}[Y]$ . Note also that with  $r_* = \mathbb{E}[Y]$ , the MSE of this approximation is the variance of  $Y$ , i.e.,

$$\mathbb{E}[(Y - r_*)^2] = \mathbb{E}[(Y - \mathbb{E}[Y])^2] = \text{Var}[Y].$$

Next, we assume that the output random variable  $Y$  is related to an input random variable  $X$  and there is a joint density function  $f_{XY}(x, y)$ . To be precise,  $X$  and  $Y$  need not have any apparent relationship, all we assume is a joint distribution  $f_{XY}(x, y)$ .

We now ask, what is the best approximation of  $Y$  by a function  $r(X)$ , i.e., what choice of the function  $r$  will minimize

$$E_{XY}[(Y - r(X))^2] = \int \int (y - r(x))^2 f_{XY}(x, y) dx dy.$$

One can minimize the above integral to find the function  $r_*$ , however, the same proof can be replicated in essence using the tower law of expectation. (For a proof based on integral minimization, see [5, pp. 263]). We re-write the total expectation using the tower law which allows us to use the result we just derived for approximating  $Y$  by a constant:

$$\begin{aligned} E_{XY}[(Y - r(X))^2] &= E_X[E_{Y|X}[(Y - r(X))^2|X]] \\ &\geq E_X[E_{Y|X}[(Y - r_*(X))^2|X]] \\ &= E_X[E_{Y|X}[(Y - E[Y|X])^2|X]] \\ &= E_X[\text{Var}[Y|X]]. \end{aligned}$$

The inner expectation is conditioned on  $X$ , therefore  $r(X)$  is a constant and by the previous result, the best  $r(X)$ , i.e.  $r_*(X)$  in such a case is given by,

$$r_*(X) = E[Y|X].$$

Therefore, the best estimate of  $Y$  given  $X$  in the MSE sense is  $E[Y|X]$ . The function  $E[Y|X]$  is called the *regression function* or the *conditional expectation function*.

**Theorem 2** (Orthogonal Decomposition). *Let  $X$  and  $Y$  be two random variables with joint distribution  $f_{XY}(x, y)$ . We can decompose  $Y$  as*

$$Y = E[Y|X] + \epsilon,$$

*such that,*

$$E[\epsilon|X] = 0,$$

*and for any function  $g(X)$*

$$E[g(X)\epsilon] = 0.$$

*any function  $g$ .*

*Proof.* The decomposition

$$Y = E[Y|X] + \epsilon,$$

only requires the existence of conditional means, which we are assuming implicitly. We have little interest in considering odd cases of distributions where means and variances are not finite. At a more philosophical level, one can observe the values of  $X$  and  $Y$  and define  $\epsilon$  by subtracting from  $Y$  the conditional expectation  $E[Y|X]$ .

The important insight of the theorem lies in the fact that the “error”  $\epsilon$  is orthogonal to the best approximation  $E[Y|X]$ . To prove this, we take conditional expectation on either side of the above decomposition. This results in

$$\begin{aligned} E[Y|X] &= E[E[Y|X]|X] + E[\epsilon|X] \\ &= E[Y|X] + E[\epsilon|X], \end{aligned}$$

which implies  $E[\epsilon|X] = 0$ .

For the second assertion, consider the tower law of iteration,

$$\begin{aligned} E[g(X)\epsilon] &= E[E[g(X)\epsilon|X]] \\ &= E[g(X)E[\epsilon|X]] \\ &= E[g(X) \cdot 0] \\ &= 0, \end{aligned}$$

where the second last equation follows from our previous result:  $E[\epsilon|X] = 0$ .  $\square$

**Corollary 1.** *Under the above decomposition  $E[\epsilon] = 0$ .*

*Proof.* Choose  $g(X) \equiv 1$  in the previous theorem.  $\square$

So far, we have just presented the conditional expectation function  $E[Y|X]$  without any assumptions on how  $X$  and  $Y$  might be related.

Linear regression assumes that the conditional expectation is *linear* in  $X$ , i.e.

$$E[Y|X] = \beta_0 + \beta_1 X,$$

Or equivalently, substituting the above into the decomposition of the previous theorem, we can say that a linear model assumes that  $X$  and  $Y$  are linearly related:

$$Y = \beta_0 + \beta_1 X + \epsilon,$$

where  $\epsilon$  is a noise term, such that:

$$E[\epsilon|X] = 0.$$

All the assumptions of the linear model are as follows:

- (1) error term  $\epsilon$  has zero conditional mean:  $E[\epsilon|X] = 0$
- (2) errors are uncorrelated:  $E[\epsilon_i \epsilon_j | X] = 0$  for  $i \neq j$ .
- (3) errors have a constant variance:  $E[\epsilon^2 | X] = \text{Var}[\epsilon | X] = \sigma^2$ .

The simplest case of linear regression is when we assume that both  $X$  and  $Y$  are one dimensional real random variables and we have a set of discrete observations,

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i,$$

such that

- (1) errors  $\epsilon_i$  have zero mean:  $E[\epsilon_i | X] = 0$
- (2) errors are uncorrelated:  $E[\epsilon_i \epsilon_j | X] = 0$  for  $i \neq j$ .
- (3) errors have constant variance:  $E[\epsilon_i^2 | X] = \text{Var}[\epsilon_i | X] = \sigma^2$ .

Note that  $E[\epsilon_i | X] = 0$  implies  $E[\epsilon_i] = 0$ , since by the tower law,

$$E[\epsilon_i] = E[E[\epsilon_i | X]] = E[0] = 0.$$

We can also show, using (1) from above and using the tower law, that the residuals  $\epsilon_i$  are uncorrelated to any function  $f(X)$ ,

$$E[f(X)\epsilon_i] = E[E[f(X)\epsilon_i | X]] = E[f(X)E[\epsilon_i | X]] = 0.$$

A key idea in motivating regression is the decomposition

$$Y = E[Y|X] + \epsilon.$$

In matrix notation, we can write the above as,

$$y = [\mathbf{1} \quad x] \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \epsilon,$$

where  $x, y$  and  $\epsilon$  are vectors in  $\mathbb{R}^n$  and  $\mathbf{1}$  is a  $n$ -vector of all ones.

$$[\mathbf{1} \quad x]^T y = [\mathbf{1} \quad x]^T [\mathbf{1} \quad x] \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + [\mathbf{1} \quad x]^T \epsilon,$$

or,

$$\begin{bmatrix} \mathbf{1}^T \\ x^T \end{bmatrix} y = \begin{bmatrix} \mathbf{1}^T \\ x^T \end{bmatrix} [\mathbf{1} \quad x] \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \mathbf{1}^T \\ x^T \end{bmatrix} \epsilon.$$

This gives,

$$\begin{bmatrix} \mathbf{1}^T y \\ x^T y \end{bmatrix} = \begin{bmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T x \\ \mathbf{1}^T x & x^T x \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \mathbf{1}^T \epsilon \\ x^T \epsilon \end{bmatrix}.$$

Let's put hats on  $\beta_0$  and  $\beta_1$ , as we are going to drop the terms involving  $\epsilon$  and apply Cramer's rule to solve the system.

$$\begin{aligned} \hat{\beta}_0 &= \frac{\begin{vmatrix} \mathbf{1}^T y & \mathbf{1}^T x \\ x^T y & x^T x \end{vmatrix}}{\begin{vmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T x \\ \mathbf{1}^T x & x^T x \end{vmatrix}} \\ &= \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - \sum x_i \sum x_i} \end{aligned}$$

$$\begin{aligned} \hat{\beta}_1 &= \frac{\begin{vmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T y \\ \mathbf{1}^T x & x^T y \end{vmatrix}}{\begin{vmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T x \\ \mathbf{1}^T x & x^T x \end{vmatrix}} \\ &= \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - \sum x_i \sum x_i} \\ &= \frac{\frac{1}{n} \sum x_i y_i - \left(\frac{1}{n} \sum x_i\right) \left(\frac{1}{n} \sum y_i\right)}{\frac{1}{n} \sum x_i^2 - \left(\frac{1}{n} \sum x_i\right) \left(\frac{1}{n} \sum x_i\right)} \\ &= \frac{\text{Cov}(x, y)}{\text{Cov}(x, x)} \\ &= \frac{\text{Cov}(x, y)}{\text{Var}(x)} \\ &= \rho \frac{\sigma_y}{\sigma_x} \end{aligned}$$



Here is the key idea to remember, the slope of the regression line  $\beta_1$  is estimated by dividing the sample covariance  $\text{Cov}(x, y)$  with the sample variance  $\text{Var}(x)$ .

Notice that the explicit solution of  $\beta_0$  above does not have a clear interpretation. However, if we solve for  $\beta_1$  first, we can find  $\beta_0$  by,

$$\beta_0 = \bar{y}_i - \beta_1 \bar{x}_i$$

## 12. STATISTICS

### 12.1. The trio: $Z$ - $\chi^2$ - $t$ .

12.1.1. *The Normal distribution.* This is just the central limit theorem. If  $X_1, X_2, \dots, X_n$  are i.i.d., with the mean and variance of each  $X_i$  being  $\mu$  and  $\sigma^2$ , respectively, then,

$$Z := \lim_{n \rightarrow \infty} \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}},$$

where  $Z$  is a standard normal.

12.1.2. *The Chi-squared distribution.* If  $Z_1, Z_2, \dots, Z_n$  are independent standard normals, then,

$$\chi_n^2 := \sum_i^n Z_i^2,$$

has a chi-square distribution with  $n$  degrees of freedom.

12.1.3. *The  $t$ -distribution.* Let  $Z$  and  $Z_1, Z_2, \dots, Z_n$ , be independent standard normals, then

$$t := \frac{Z}{\sqrt{\frac{\sum_{i=1}^n Z_i^2}{n}}} = \frac{Z}{\sqrt{\frac{\chi_n^2}{n}}},$$

has a  $t$ -distribution.

12.2. **Hypothesis testing.** Type-I error is the error of rejecting a null hypothesis when it is in fact true. The probability of this error is usually denoted by  $\alpha$ . A typical value of  $\alpha$  used in practice is .05.

### 12.3. The Law of Total Expectation.

$$E(X) = E(E(X|Y)),$$

or to be more explicit regarding the probability densities used for each expectation:

$$E(X) = E_Y(E_{X|Y}(X|Y)).$$

#### 12.4. The Law of Total Variance.

$$\text{Var}(X) = \text{Var}(\mathbb{E}(X|Z)) + \mathbb{E}(\text{Var}(X|Z))$$

*Proof.* By definition,

$$\text{Var}(X|Z) = \mathbb{E}(X^2|Z) - \mathbb{E}(X|Z)^2.$$

Taking expectation w.r.t  $Z$ , we get,

$$\begin{aligned} \mathbb{E}(\text{Var}(X|Z)) &= \mathbb{E}(\mathbb{E}(X^2|Z)) - \mathbb{E}(\mathbb{E}(X|Z)^2) \\ &= \mathbb{E}(X^2) - \mathbb{E}(\mathbb{E}(X|Z)^2) \\ &= \mathbb{E}(X^2) - \mathbb{E}(X)^2 - (\mathbb{E}(\mathbb{E}(X|Z)^2) - \mathbb{E}(X)^2) \\ &= \text{Var}(X) - (\mathbb{E}(\mathbb{E}(X|Z)^2) - \mathbb{E}(\mathbb{E}(X|Z))^2) \\ &= \text{Var}(X) - \text{Var}(\mathbb{E}(X|Z)). \end{aligned}$$

□

### 13. MATRICES

#### 13.1. Rank Theorem.

**Theorem 3.** *Given a matrix  $A$ , the row rank of  $A$  is the same as the column rank of  $A$ .*

*Proof.* Convert  $A$  to its row reduced echelon form  $R$ . This is done through elementary matrix operations on  $A$ ,

$$\begin{aligned} E_n \dots E_2 E_1 A &= UA = R, \\ A &= U^{-1}R. \end{aligned}$$

Therefore, rows of  $R$  can be obtained as a linear combination of the rows of  $A$  and vice-versa, hence,

$$\text{Row}(A) = \text{Row}(R).$$

The columns of  $A$  and  $R$  are related as follows,

$$\begin{aligned} UA &= U [C_1 \ C_2 \ \dots \ C_n] \\ &= [UC_1 \ UC_2 \ \dots \ UC_n] \\ &= [C_1^R \ C_2^R \ \dots \ C_n^R] \\ &= R. \end{aligned}$$

If the row rank of  $A$  is  $r$ , then the matrix  $R$  has  $r$  rows with leading 1s. The key insight is that the columns of  $R$  with leading 1s form a basis of  $\text{Col}(R)$ .

We claim that the columns of  $A$  for which  $R$  has a leading 1, form a basis of  $\text{Col}(A)$ . That is, the subset of columns of  $A$ ,

$$\phi = \{C_{j_1}, C_{j_2}, \dots, C_{j_r}\},$$

where  $j_1, j_2, \dots, j_r$  are the indices of columns where  $R$  has a leading 1, is a basis of  $\text{Col}(A)$ .

To prove the linear independence of  $\phi$ , let

$$a_1 C_{j_1} + a_2 C_{j_2} + \dots + a_r C_{j_r} = 0,$$

then

$$\begin{aligned} a_1 U C_{j_1} + a_2 U C_{j_2} + \dots + a_r U C_{j_r} &= 0, \\ a_1 C_{j_1}^R + a_2 C_{j_2}^R + \dots + a_r C_{j_r}^R &= 0. \end{aligned}$$

Since the columns  $\{C_{j_1}^R, C_{j_2}^R, \dots, C_{j_r}^R\}$  of  $R$  are a basis of  $\text{Col}(R)$ , the last equation implies that

$$a_1 = a_2 = \dots = a_r = 0.$$

Therefore,  $\phi$  is a linearly independent set.

To prove that  $\phi$  spans  $\text{Col}(A)$ , consider  $b \in \text{Col}(A)$ . Then, there exists  $x$  such that,

$$b = Ax = U^{-1}Rx.$$

Since  $Rx \in \text{Col}(R)$  and  $\{C_{j_1}^R, C_{j_2}^R, \dots, C_{j_r}^R\}$  is a basis of  $\text{Col}(R)$ , we can write the last equation as,

$$\begin{aligned} b &= U^{-1}(t_{j_1} C_{j_1}^R + t_{j_2} C_{j_2}^R \dots + t_{j_r} C_{j_r}^R) \\ &= t_{j_1} C_{j_1} + t_{j_2} C_{j_2} \dots + t_{j_r} C_{j_r}, \end{aligned}$$

and hence  $b \in \text{Span}(\phi)$ .

This proves that  $\phi$  is a basis of  $\text{Col}(A)$  and therefore,

$$\begin{aligned} \dim(\text{Col}(A)) &= r = \dim(\text{Row}(A)) \\ \text{Row rank}(A) &= r = \text{Col Rank}(A). \end{aligned}$$

□

**Lemma 1.** *A and  $A^T$  have the same rank.*

**13.2. Left Right Inverses.** If a square matrix  $A$  has a left inverse, prove that  $A$  is invertible, i.e.,  $A$  has a right inverse that matches the left inverse.

*Proof.* Let  $B_L$  be the left inverse of  $A$ . Then, by definition,

$$B_L A = I.$$

This implies that the columns of  $A$  are linearly independent, since  $Ax = 0 \implies x = 0$ :

$$\begin{aligned} Ax &= 0 \\ B_L Ax &= 0 \\ Ix &= 0 \\ x &= 0. \end{aligned}$$

Therefore,  $A$  has full column rank. Let  $n$  be the number of columns of  $A$ , then the columns of  $A$  are a basis of  $\mathbb{R}^n$ ,

- (1) Because  $A$  is square, each column of  $A$  is in  $\mathbb{R}^n$ . (This is the only place where we use the fact that  $A$  is square).
- (2) Columns of  $A$  form a set of  $n$  linearly independent vectors in  $\mathbb{R}^n$  and therefore must span  $\mathbb{R}^n$ .

This proves that the columns of  $A$  can be linearly combined to form any given right hand side, that is, for a right hand side  $y$ , we can find a solution  $x$  such  $Ax = y$ . Using the left inverse, one possible solution is given by  $x = B_L y$ . The solution is unique, for if  $x'$  is another solution, then,

$$\begin{aligned} Ax' &= y = Ax \\ x' &= B_L y = x. \end{aligned}$$

Choosing,  $y = e_i$  for  $i = 1, 2, \dots, n$ , the above implies the existence of a unique matrix  $B_R$  where,

$$\begin{aligned} B_R &= [x_1 \quad x_2 \quad \dots \quad x_n] \\ &= [B_L e_1 \quad B_L e_2 \quad \dots \quad B_L e_n] \\ &= B_L, \end{aligned}$$

such that,

$$AB_R = [e_1 \quad e_2 \quad \dots \quad e_n] = I.$$

This proves the existence of a right inverse  $B_R$  and shows that  $B_R = B_L$ .

For another way to show that  $B_L = B_R$ , compute the product  $B_L AB_R$  in the following two ways,

$$\begin{aligned} B_L(AB_R) &= B_L I = B_L, \\ (B_L A)B_R &= I B_R = B_R, \end{aligned}$$

which proves  $B_L = B_R$ . □

Therefore,  $A$  is invertible since  $B := B_L = B_R$  is the unique inverse of  $A$ .

**13.3. LU Factorization.** Given a square  $n \times n$  matrix  $A$ , we can find the  $LU$  decomposition of  $A$  such that

$$(13.1) \quad PA = LU,$$

where  $P$  is a permutation matrix,  $L$  is a lower triangular and  $U$  is an upper triangular.

The algorithm to compute the matrices  $L$  and  $U$ , given  $A$  is essentially the row reduced echelon form computation of  $A$ . We want to apply a sequence of matrices  $L_i$  such that  $A$  gets transformed into  $U$ .

$$(13.2) \quad L_1 L_2 \dots L_{n-1} PA = U,$$

where each  $L_i$  is a lower triangular matrix which introduces zeros in the  $i^{th}$  column below the diagonal.

**13.4. Unitarily Diagonalizable Matrices.** A unitarily diagonalizable matrix is called a *Normal* matrix.

$$A \text{ is normal} \Leftrightarrow A = Q\Lambda Q^*,$$

where  $Q$  is unitary i.e.,  $Q^*Q = I = QQ^*$  and  $\Lambda$  is diagonal.

A normal matrix commutes with its complex conjugate, and this property also characterizes a normal matrix.

$$A \text{ is normal} \Leftrightarrow A^*A = AA^*.$$

Orthogonally Diagonalizable Matrices	
Complex	Real
Unitary	Orthogonal
Hermitian	Symmetric
Skew-Hermitian	Skew-Symmetric

**Real Quadratic Form.** Let  $A$  be a square matrix in  $\mathbb{C}^{n \times n}$ . If,

$$x^*Ax \in \mathbb{R}, \quad \forall x \in \mathbb{C}^n,$$

then  $A$  is Hermitian.

*Proof.* Let

$$B = \frac{A + A^*}{2}, C = \frac{A - A^*}{2i},$$

then  $A = B + iC$  and both  $B$  and  $C$  are Hermitian. Since  $x^*Ax$  is real, we have,

$$\begin{aligned} x^*Ax &= (x^*Ax)^* \\ x^*(B + iC)x &= x^*(B^* + (iC)^*)x \\ x^*(B + iC)x &= x^*(B - iC)x \\ x^*Cx &= -x^*Cx \\ x^*Cx &= 0. \end{aligned}$$

We now show that if  $x^*Cx = 0$  for all  $x \in \mathbb{C}^n$ , then  $C$  is the zero matrix.

For any vectors  $x, y$ , we have,

$$\begin{aligned} (x + y)^*C(x + y) &= 0 \\ x^*Cx + y^*Cy + x^*Cy + y^*Cx &= 0 \\ 0 + 0 + x^*Cy + y^*Cx &= 0 \\ x^*Cy + (x^*Cy)^* &= 0 \\ \operatorname{Im}(x^*Cy) &= 0. \end{aligned}$$

Replace  $y$  by  $iy$  in the above and we get

$$\begin{aligned} \operatorname{Im}(ix^*Cy) &= 0 \\ \operatorname{Re}(x^*Cy) &= 0. \end{aligned}$$

Therefore,  $x^*Cy = 0$  for any  $x, y$ . Now choose  $x, y$  to be the unit vectors along directions  $i, j$  to obtain  $e_i^*Ce_j = c_{ij} = 0$  which implies that  $C$  is the zero matrix and hence  $A = A^*$ .

□

**13.5. Cholesky Decomposition in Finance.** The fundamental use of Cholesky decomposition in finance is to generate a vector of correlated samples from a vector of uncorrelated samples, where the correlations are specified by the correlation matrix. This correlation matrix, say  $A$  is by construction real, symmetric, and positive semi-definite.

13.5.1. *Example:* Suppose we have two independent zero mean random variables  $Z_1$  and  $Z_2$  with variances  $\sigma_1$  and  $\sigma_2$ , respectively. We wish to construct random variables  $X_1$  and  $X_2$ , again with variances  $\sigma_1$  and  $\sigma_2$  respectively, such that the correlation between  $X_1$  and  $X_2$  is  $\rho$ .

The correlation matrix in this simple  $2 \times 2$  case is

$$(13.3) \quad A = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}.$$

The Cholesky decomposition is very simple (think RREF):

$$(13.4) \quad A = \begin{bmatrix} 1 & 0 \\ \rho & 1 \end{bmatrix} \begin{bmatrix} 1 & \rho \\ 0 & 1 - \rho^2 \end{bmatrix}.$$

We can write this in the form where diagonal of  $L$  and  $U$  both are replaced by the squar root of their products, and we get

$$(13.5) \quad A = \begin{bmatrix} 1 & 0 \\ \rho & \sqrt{1 - \rho^2} \end{bmatrix} \begin{bmatrix} 1 & \rho \\ 0 & \sqrt{1 - \rho^2} \end{bmatrix}.$$

It can be shown, that applying  $L$  on the vector  $z$  gives us the required correlated vector  $x$ :

$$(13.6) \quad \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \rho & \sqrt{1 - \rho^2} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}.$$

**Perturbation of the Identity Matrix.**

$$(I + uv^T)^{-1} = I - \frac{uv^T}{1 + v^T u}.$$

$$(I + UV)^{-1} = I - U(I + VU)^{-1}V.$$

**Condition number.** The sensitivity of fractional change in output to fractional change in input is called the *condition number* of a problem.

Suppose we have the problem of solving  $Ax = b$  for a square, non-singular matrix  $A$ . There are a number of ways one can define the output and input variables. One useful and practical choice is to consider  $A$  as the input,  $x = A^{-1}b$  as the output, and assume that  $b$  remains constant as the input and output change. The problem is to find the mapping  $f : A \rightarrow x = A^{-1}b$ . By definition, the condition number of this problem is given by,

$$(13.7) \quad \kappa(f) = (\|\Delta x\|/\|x\|) / (\|\Delta A\|/\|A\|).$$

Since  $b$  is constant as input and output change, we get,

$$(A + \Delta A)(x + \Delta x) = b.$$

Using  $Ax = b$  and ignoring  $\Delta A \Delta x$  to get,

$$A\Delta x = -\Delta Ax,$$

or

$$\Delta x = -A^{-1}\Delta Ax,$$

which implies

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A^{-1}\| \|\Delta A\|.$$

This together with (13.7) gives the bound

$$\kappa(f) \leq \|A\| \|A^{-1}\|.$$

The number  $\kappa(A) := \|A\| \|A^{-1}\|$  appears in so many other contexts and problems that it is called the condition number of  $A$ . In the 2-norm,  $\kappa(A)$  is the ratio of the largest to the smallest singular value of  $A$ .

**Projection Matrices.** A square matrix  $P$  such that  $P^2 = P$  is called a projection matrix.  $P$  projects any vector on its column space.  $P$  is an idempotent matrix by definition and vice-versa an idempotent matrix is a projection matrix projecting on its own column space.

*Trace-Eigenvalues-Rank of Projection Matrices.* An eigenvalue of a projection matrix is either 1 or 0.

*Proof.* Let  $P$  be a projection matrix and  $\lambda$  be an eigenvalue of  $P$  with eigenvector  $x$ . Then,

$$(13.8) \quad \lambda x = Px = P^2x = P(Px) = P(\lambda x) = \lambda Px = \lambda^2 x.$$

Therefore,  $\lambda = \lambda^2$ , which implies that either  $\lambda = 1$  or  $\lambda = 0$ . □

Since trace equals the sum of eigenvalues, for projection matrices, the trace equals the rank.

*Proof.* Let  $P$  be an  $n \times n$  projection matrix of rank  $r$ . Then  $P$  has  $n - r$  linearly independent vectors in its null space, i.e., the eigenvalue 0 has a multiplicity of  $n - r$ . Each of the remaining  $r$  eigenvectors of  $P$  must have eigenvalue 1 since the eigenvalues of  $P$  are contained in the set  $\{0, 1\}$ . Therefore,

$$\begin{aligned} \text{tr}(P) &= \sum_{i=1}^n \lambda_i \\ &= 1 \times r + (n - r) \times 0 \\ &= r. \end{aligned}$$

□

The only full rank projection matrix is the identity matrix.

*Proof.*  $P$  is full rank and square, therefore  $P^{-1}$  exists. Since  $P^2 = P$ , multiplying both sides by  $P^{-1}$  gives  $P = I$ .  $\square$

A typical example of a projection matrix is the rank-1 matrix,

$$\frac{aa^T}{a^T a},$$

where  $a$  is a vector in  $\mathbb{R}^n$ . The projection of a vector  $b$  on  $a$  is

$$Proj_a(b) = \frac{aa^T}{a^T a} b = \frac{a^T b}{a^T a} a.$$

The eigenvectors of  $\frac{aa^T}{a^T a}$  are:

- (1) All vectors orthogonal to  $a$  (each with eigenvalue 0).
- (2) The vector  $a$  (with eigenvalue 1).

In general, for a full column rank matrix  $A$ , the matrix

$$A(A^* A)^{-1} A^*$$

will project on the columns of  $A$ . In particular, the projection of a vector  $b$  on the column space of  $A$  is given by

$$Proj_{Col(A)}(b) = A(A^* A)^{-1} A^* b.$$

Given  $b$  and  $A$ , the least squares problem is to find an  $x$  such that  $Ax - b$  has the smallest 2-norm. By definition, projecting  $b$  on the column space of  $A$  makes the projection error  $Ax - b$  orthogonal to the column space of  $A$ , hence minimising the 2-norm of the error. Therefore, the least squares problem boils down to a projection,

$$(13.9) \quad Ax = Proj_{Col(A)}(b) = A(A^* A)^{-1} A^* b.$$

Since  $A$  has full column rank, the left inverse of  $A$  exists and we immediately get

$$x = (A^* A)^{-1} A^* b.$$

Equivalently, we can re-write (13.9) as

$$A(x - (A^* A)^{-1} A^* b) = 0,$$

and since  $A$  has full column rank, its null space is trivial, therefore,

$$x = (A^* A)^{-1} A^* b.$$

If  $Q$  is a unitary matrix, then the associated projection matrix is  $QQ^*$ . The decomposition  $QR = A$ , together with the projection matrix  $QQ^*$  leads to another algorithm for the least squares problem. We have,

$$\begin{aligned} Ax &= Proj_{Col(A)}(b) \\ QRx &= QQ^* b \\ Rx &= Q^* b. \end{aligned}$$

The last displayed system is upper triangular and can be solved recursively starting from the bottom. The  $QR$  decomposition serves as one of the best numerical methods for solving the least squares problem.



**Variance Covariance Matrices.** Let  $X$  be a random vector such that

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix},$$

such that the variance covariance matrix of  $X$  is given by  $V$  and let  $a$  and  $b$  be vectors in  $\mathbb{R}^n$ , then

$$\text{Cov}(a^T X, b^T X) = a^T V b = b^T V a.$$

If  $A$  is a matrix in  $\mathbb{R}^n$ , then the covariance of  $AX$  is  $AVAT$ .

#### 14. USING CHEBYSHEV METHODS

- (1) write a basic core
- (2) must have a root finder

Where can it be used?

- (1) Corr Bump: It's a 1d root finding problem

#### INEQUALITIES

**Bessel's Inequality.** Projection of a vector is at most as big as the vector.

*Proof.* Let  $x_j$  be an orthonormal set in a Hilbert space  $\mathcal{H}$ . For any  $x \in \mathcal{H}$ , we have,

$$0 \leq \left\| x - \sum_j (x, x_j) x_j \right\|^2 = \|x\|^2 - \sum_j |(x, x_j)|^2.$$

□

**Cauchy–Schwarz Inequality.** The cosine of the angle between two vectors is between  $-1$  and  $1$ .

$$|(x, y)|^2 \leq \|x\|^2 \|y\|^2.$$

*Proof.* If  $y = 0$ , the inequality is trivial. If  $y \neq 0$ , define  $y_0 = \frac{y}{\|y\|}$ . Project  $x$  on  $y_0$  and apply Bessel's inequality, which gives,

$$0 \leq \|x - (x, y_0) y_0\|^2 = \|x\|^2 - |(x, y_0)|^2.$$

□

## QUESTIONS

- (1) Six vectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_6$  are chosen to be either  $[1, 1]^T$  or  $[3, 2]^T$  with equal probability, with each choice made independently. What is the probability that the sum  $\mathbf{v}_1 + \mathbf{v}_2 + \dots + \mathbf{v}_6$  is equal to the vector  $[10, 8]^T$  [4] ?

Let's first find the number of ways we can sum to  $[10, 8]^T$ . Let  $n$  be the number of times  $[1, 1]^T$  is chosen. Then  $[3, 2]^T$  is chosen  $6 - n$  times. Setting  $n[1, 1]^T + (6 - n)[3, 2]^T = [10, 8]$  gives  $n = 4$  as the only choice.

The probability of choosing the vector  $[1, 1]^T$  4 times is the same as the probability of 4 heads in 6 independent tosses of a fair coin:

$$P(k = 4) = \binom{6}{4} (0.5)^4 (0.5)^{6-4} = \frac{15}{64}.$$

- (2) We have two matrices  $A$  and  $B$ , both having the same number of rows. How can we check if the column space of  $A$  and the column space of  $B$  are equal?

Find the row reduced echelon form (RREF) of both  $A^T$  and  $B^T$  and discard the zero rows. The column spaces of  $A$  and  $B$  are equal iff

$$\begin{aligned} \text{RREF}(A^T) &= \text{RREF}(B^T) \\ U_n U_{n-1} \dots U_2 U_1 A^T &= V_m V_{m-1} \dots V_2 V_1 B^T \\ U A^T &= V B^T, \end{aligned}$$

where  $U_i, V_j, U, V$  are invertible.

- (3) Evaluate  $\int_0^1 \frac{x-1}{\ln x} dx$ .

Use Feynman's trick. Let

$$I(\alpha) = \int_0^1 \frac{x^\alpha - 1}{\ln x} dx, \quad \alpha \in (0, 1).$$

Note that  $I(0) = 0$ . Differentiate  $I$  w.r.t  $\alpha$ , we get

$$\begin{aligned} I'(\alpha) &= \int_0^1 x^\alpha dx, \quad \alpha \in (0, 1) \\ &= \frac{1}{\alpha + 1}. \end{aligned}$$

Integrating the last result again, we get

$$\begin{aligned} I(\alpha) &= \ln(\alpha + 1) + I(0) \\ &= \ln(\alpha + 1). \end{aligned}$$

The last equation implies,

$$\lim_{\alpha \rightarrow 1} I(\alpha) = \ln(2).$$

- (4) Suppose  $X$  and  $Y$  are positive random variables. What is the relationship between  $P(X + Y > 4)$  and  $P(X > 2)P(Y > 2)$ ? Is one greater, smaller than the other or no such relationship in general?

Use  $X = Y$  and  $X = -Y$  and see... but both  $X, Y$  have to be positive, so that doesn't really work or does it?

- (5) We are given an array of integers starting from 0 and ending at 100. We start from 0 and start flipping a fair coin. If a heads come down, we advance by 1 and if a tail comes down we advance by 2. For example, if head comes down, we move to index 1, otherwise, we skip 1 and move to index 2. We continue flipping the coin until we reach hundred or beyond.

The question is, which index in the array has the largest probability of being visited?

- (6) Let  $X, Y$  and  $Z$  be random variables. If the correlations between them satisfy the following relationship,

$$\rho_{XY} = \rho = \rho_{YZ}$$

what can you say about the correlation  $\rho_{XZ}$ ?

Thinking in terms of vectors, let  $\alpha$  be the angle between  $X$  and  $Y$ , which must be the same as the angle between  $Y$  and  $Z$ . Then

$$\cos \alpha = \rho,$$

[TODO]: This can be done through looking at the positive semi-definite property of the correlation matrix as well.

[TODO]: Draw a picture of  $X, Y$  and  $Z$ . If  $X$  and  $Y$  are orthogonal, i.e.,  $\rho = 0$ , then  $Y$  must be normal to the plane formed by  $Z$  and  $X$ . In this case,  $\rho_{XZ}$  can be any number between  $-1$  and  $1$ .

If  $X, Y$  are parallel, i.e.  $\rho = 1$ , and  $Y, Z$  are parallel, then  $X, Z$  must be parallel as well. (Draw three vectors pointing in the same direction). Therefore,  $\rho_{XZ} = 1$ .

If  $X, Y$  are anti-parallel, i.e.  $\rho = -1$ , and  $Y, Z$  are anti-parallel, then  $X, Z$  must be parallel and therefore,  $\rho_{XZ} = 1$ .

For the general case, if  $X$  and  $Y$  make an angle  $\alpha$  between each other and  $Y$  and  $Z$  make the same angle  $\alpha$  between each other, then the angle between  $X$  and  $Z$  can be between  $0$  and  $2\alpha$ . [TODO: Draw cones]

Therefore, the range of correlation between  $X$  and  $Z$  must satisfy

$$\begin{aligned} \cos(2\alpha) &\leq \rho_{XZ} \leq \cos(0), \\ \implies \cos(2\cos^{-1}(\rho)) &\leq \rho_{XZ} \leq 1. \end{aligned}$$

- (7) Suppose you are trading and you define your *hit ratio* as the number of trades you have won divided by the total number of trades you participated in. You checked your hit ratio in the morning and it was less than  $0.9$ . You checked later in the day and it was larger than  $0.9$ . Is it the case that at some point during the day, the hit ratio was exactly  $0.9$ ?

Yes. Let  $k_1$  be the number of trades lost out of a total of  $n_1$  trades. Then, in the morning you had

$$\frac{n_1 - k_1}{n_1} < 0.9,$$

which implies

$$n_1 < 10k_1,$$

and similarly, later in the day you had

$$\frac{n_2 - k_2}{n_2} > 0.9,$$

$$n_2 > 10k_2.$$

where  $n_2 \geq n_1$  and  $k_2 \geq k_1$ .

The variables  $n_1 \leq n_2$  being discrete must be equal to  $10k$  at some intermediate time at which point your hit ratio must be 0.9.

There is nothing special about 0.9, the same can be concluded for any hit ratio  $r$  as long as  $r$  has a specific form. Let

$$\frac{n - k}{n} = r,$$

which implies,

$$n = \frac{1}{1-r}k.$$

Therefore, as long as  $\frac{1}{1-r}$  is an integer, the problem structure remains the same.

Let  $\frac{1}{1-r} = j$  where  $j \in \mathbb{Z}^+$ , then the permissible values of  $r$  are:

$$r = 1 - \frac{1}{j}, \quad j \in \mathbb{Z}^+.$$

- (8) Given a data series,  $x_1, x_2, \dots, x_n$ , for a large  $n$ , how do you test if the data is normal? Given another data series,  $y_1, y_2, \dots, y_n$ , how do you compare the normality of the two data series with each other, i.e. which data series is *more* normal than the other?
- (9) You have a jury of three judges. Two of the judges make the correct decision, each with probability  $p$ , while the third judge tosses a fair coin to make the decision. All judges think/toss coins independently. What is the probability that the jury makes the right decision?

Let  $C$  be the event that the jury makes the correction decision. Then, conditioning on whether the judge with the coin makes the correct decision ( $j_3$ ), we have

$$P(C) = P(C|j_3)P(j_3) + P(C|\neg j_3)P(\neg j_3).$$

This boils down to

$$P(C) = \frac{1}{2} (1 - (1-p)^2) + \frac{1}{2} p^2 = p.$$

Thanks to the coin tossing, a three person jury has effectively been reduced to a one person jury.

- (10) What is the relationship of  $\text{Var}(\text{E}(X | Y))$  with  $\text{E}(\text{var}(X | Y))$ ?

The law of total variance:

$$\text{Var}(Y) = \text{E}[\text{Var}(Y | X)] + \text{Var}(\text{E}[Y | X]).$$

- (11) Find the value of  $\alpha$  which minimizes,

$$\sum_i^n |\alpha w_i - n_i|,$$

where  $w_i$  and  $n_i$  are positive real numbers.

Notice that  $\alpha w_i - n_i$  has a zero at  $\frac{n_i}{w_i}$ . Define

$$a := \min \left\{ \frac{n_i}{w_i} : i = 1, 2, \dots, n \right\},$$

and

$$b := \max \left\{ \frac{n_i}{w_i} : i = 1, 2, \dots, n \right\}.$$

Now apply a bisection on the interval  $[a, b]$ .

- (12) Estimate the derivative of  $x^x$  at  $x = 2$ .

We have:

$$\begin{aligned} \frac{d}{dx} x^x &= \frac{d}{dx} e^{x \ln x} \\ &= e^{x \ln x} \frac{d}{dx} (x \ln x) \\ &= x^x \frac{d}{dx} (\ln x + 1). \end{aligned}$$

At  $x = 2$ , the last expression is roughly equal to  $2^2(.7 + 1) = 6.8$ .

- (13) What are the eigenvalues of an  $n \times n$  matrix all of whose entries are a constant?

The rank of this matrix is 1, so 0 is an eigenvalue with multiplicity  $n - 1$ . The only non-zero eigenvalue is  $n$  and the corresponding eigenvector is a vector of all ones.

- (14) Let  $X$  be a random vector in  $\mathbb{R}^n$  whose  $n \times n$  variance-covariance matrix is  $\Sigma$ . What is the variance-covariance matrix of  $AX$  where  $A$  is an  $n \times n$  matrix?

$$\text{Var}(AX) = A\Sigma A^T$$

An important example of the above is when  $X = Z$  i.e., a random vector of i.i.d standard normals. In that case  $\Sigma = I$  and the variance-covariance matrix of  $AZ$  is  $AA^T$ .

- (15) You roll a hundred sided die and note the number that appears. How many rolls you need on average to see the same number again?

The answer is 100. You roll the first time and get some number, say 12. Your probability of getting a 12 on any following roll is  $1/100$ . You therefore need 100 more rolls on average to get the same number again.

- (16) Roll 3 standard six-sided dice together. What is the probability that the max is less than or equal to 3.

Max less than or equal to 3 is the same as each of the three dice independently showing a number less than or equal to 3. Therefore, the probability is  $(1/2)(1/2)(1/2) = 1/8$ . (We have used the fact that the probability of a six-sided dice showing a number less than or equal to 3 is  $1/2$ .)

- (17) Let  $X$  be a discrete random variable which takes values in  $\{1, 2, 3\}$  with equal probability. What is the standard deviation of  $X$ ? If we take three independent samples  $X_1, X_2$  and  $X_3$ , what is the standard deviation of the sum  $S = X_1 + X_2 + X_3$ ?

By calculation,  $\sigma_X = \sqrt{2/3}$  and  $\sigma_S = \sqrt{2}$

- (18) Let  $X$  and  $Y$  be two random variables with  $\sigma_X = 2$ ,  $\sigma_Y = 3$ . If  $\sigma_{X+Y} = 5$ , then what is the correlation  $\rho_{XY}$ ?

We use the formula:

$$\sigma_{X+Y}^2 = \sigma_X^2 + \sigma_Y^2 + 2\rho_{XY}\sigma_X\sigma_Y,$$

which implies

$$2\sigma_X\sigma_Y = 2\rho_{XY}\sigma_X\sigma_Y$$

i.e.

$$\rho_{XY} = 1.$$

- (19) Generic version of the above question. If  $\sigma_{X+Y} = \sigma_X + \sigma_Y$ , then what is the value of  $\rho_{XY}$ ?

We have:

$$(\sigma_X + \sigma_Y)^2 = \sigma_{X+Y}^2 = \sigma_X^2 + \sigma_Y^2 + 2\rho_{XY}\sigma_X\sigma_Y,$$

where the first equality is given and the second is by definition. This simplifies to

$$2\sigma_X\sigma_Y = \rho_{XY}2\sigma_X\sigma_Y,$$

which implies  $\rho_{XY} = 1$ .

**Intuition:** The standard deviation does not add linearly for uncorrelated variables. If correlation between two variables is zero, then the *variance* of the sum is indeed additive but we have the basic identity for  $a, b > 0$ :

$$\sqrt{a+b} < \sqrt{a} + \sqrt{b},$$

from which it follows that the standard deviation of the sum would be strictly less than the sum of individual standard deviations. Given that the

standard deviations added up linearly in our case gives us a hint that the variables must be highly positively correlated.

- (20) Approximate  $\log_7 1250$
- (21) Approximate  $4^{3.6}$
- (22) the mean and variance of a random variable  $X$  is 10. What is  $E[X^2]$ ?
- (23) You pick three uniform i.i.d samples from  $[0, 1]$ . What the expected value of the product of the maximum and the minimum.
- (24) Consider the following game with a six-sided standard dice. If you roll 1, 2 or 3, you roll again. If you roll 4 or 5, you get the number of dollars equal to the number of rolls you have rolled so far, including the last roll. If you roll a 6, the game ends and you get nothing. On average, how many dollars are you expected to make by playing this game?

Let's first solve an easier game: in the modified game, you can roll again if you get 1, 2 or 3 as before but if you roll 4, 5 or 6 on your  $k^{th}$  roll, you get  $k$  dollars. In this case, you can map the game to a fair coin, since the probability of getting a value in the set  $\{1, 2, 3\}$  or in the set  $\{4, 5, 6\}$  is each  $1/2$ . The expected number of throws to get an outcome in the set  $\{4, 5, 6\}$  is therefore  $1/(1/2) = 2$  and on average you make \$2 playing this game.

Now, we come back to the original question. One thing to notice is that on average, the original game must pay less than the modified game above because in the original game you either get paid nothing or you get paid the same amount that the modified game would pay. Since the good outcomes are  $2/3$  times the outcomes of the modified game, the expected value of the original game is also  $2/3$  times the expected value of the modified game. The final answer is  $(2/3)(2) = 4/3$ .

We can also solve this problem in the classical mathematical way.

Let  $X$  be the random variable that you win  $k$  dollars by rolling a 4 or 5 at the  $k^{th}$  roll. Then  $P(X = k)$  is given by the expression,

$$P(X = k) = (1/2)^{k-1}(2/6),$$

i.e. for the first  $(k - 1)$  rolls you got 1, 2 or 3, each independently with probability  $1/2$  and then on the  $k^{th}$  go, you rolled a 4 or a 5 with probability  $2/6$ .

The expected value of  $X$  is given by:

$$\begin{aligned} E[X] &= \sum_{k=1}^{\infty} kP(X = k) \\ &= (1/3) \sum_{k=1}^{\infty} k(1/2)^{k-1} \end{aligned}$$

- (25) What are the eigenvalues of an  $n \times n$  matrix all of whose off diagonal entries are a constant and the diagonal has ones?

Let  $A$  be the matrix in question and  $c$  be the off diagonal constant. We can write  $A$  as

$$(14.1) \quad A = cW + (1 - c)I,$$

where  $W$  is a matrix of all ones and  $I$  is the identity matrix.

The eigenvalue problem can now be transformed as

$$(14.2) \quad A - \lambda I = 0 \iff W - \left( \frac{\lambda + c - 1}{c} \right) I = W - \omega I = 0,$$

where  $\omega = \frac{\lambda + c - 1}{c}$ .

Now, the eigenvalues of  $W$  are 0 and  $n$ , where  $n$  is the dimension of  $W$ . This is true since  $W$  is of rank 1 which implies zero is an eigenvalue of multiplicity  $n - 1$ . Also, a vector of all ones is an eigenvector of  $W$  with  $n$  as the eigenvalue.

Therefore,

$$(14.3) \quad \omega = 0 \iff \frac{\lambda + c - 1}{c} = 0 \implies \lambda = 1 - c.$$

And

$$(14.4) \quad \omega = n \iff \frac{\lambda + c - 1}{c} = n \implies \lambda = 1 + (n - 1)c.$$

- (26) A first order Taylor approximation of a convex function always underestimates the function:

$$(14.5) \quad e^x - (1 + x) > 0, \quad x \in \mathbb{R}.$$

Can we prove this? [TODO]

- (27) Let  $X$  be a standard normal variable,  $a \in \mathbb{R}$  and  $p = P(X \leq a)$ . Define a new random variable  $Y = -X$ . What is  $P(Y \leq a)$ ?

$$(14.6) \quad P(Y \leq a) = P(-X \leq a)$$

$$(14.7) \quad = P(X \geq -a)$$

$$(14.8) \quad P(Y \leq a) = P(X \leq a) = p.$$

In the last step, we have used the symmetry property of the cumulative distribution function of a standard normal variable.

Intuitively, the result makes sense<sup>2</sup>. A normal random variable is symmetric around the origin. Its distribution does not change when we reflect it across the origin.

- (28) Let us consider the Indicator function  $I(X > 30)$ , where  $X$  is a standard normal variable. Suppose you want to naively compute the expected value of  $I(X > 30)$ . Note that this number is strictly greater than 0, since

$$E[I(X > 30)] = P(X > 30) = \phi(-30) > 0.$$

---

<sup>2</sup>It does not have to, some results in probability are counter intuitive to most people. Also, what some people think intuitive might not be intuitive to others. Intuition is subjective.



However, if we apply the *standard* approach of approximating the expected value via the discrete sampling

$$E[I(X > 30)] \simeq \frac{1}{N} \sum_{i=1}^N I(X_i > 30),$$

what is the expected value of  $N$  to make the sum greater than zero?

*Answer:* Convert the problem to a geometric random variable, with the success probability  $p = \phi(-30)$ . Therefore  $E[N] = 1/\phi(-30)$ .

(29) What is  $\frac{1}{\sqrt{2\pi}}$  approximately?

*Answer:* 0.4.

(30) Let us suppose that  $X$  is a standard normal. What is  $E[X|X > 0]$ .

*Answer:*  $\sqrt{\frac{2}{\pi}} \simeq 0.8$ .

**Polya's Urn:** An urn has  $n_b = 95$  black balls and  $n_w = 5$  white balls. Therefore, the ratio  $r_b$  of black balls to the total number of balls in the urn is

$$(14.9) \quad r_b = \frac{n_b}{n_b + n_w} = 0.95.$$

You start a game, where you pick a ball randomly and if the ball is black, you add some more black balls and if it's white you add some more white balls. Suppose you have done this a few million times and you ask now, what is the expected value of the ratio  $r_b$  now?

*Answer:* .95. The ratio is a martingale with respect to the stochastic process. Let  $r_{bk}$  be the ratio at the  $k^{th}$  turn. At the  $(k+1)^{st}$  turn, we have

(14.10)

$$(14.11) \quad \begin{aligned} E[r_{b(k+1)}] &= \left( \frac{n_{bk}}{n_{bk} + n_{wk}} \right) \frac{n_{bk} + f(k)}{n_{bk} + n_{wk} + f(k)} + \left( \frac{n_{wk}}{n_{bk} + n_{wk}} \right) \frac{n_{bk}}{n_{bk} + n_{wk} + f(k)} \\ &= \left( \frac{n_{bk}}{n_{bk} + n_{wk}} \right) \frac{n_{bk} + n_{wk} + f(k)}{n_{bk} + n_{wk} + f(k)} \end{aligned}$$

$$(14.12) \quad = \frac{n_{bk}}{n_{bk} + n_{wk}}$$

$$(14.13) \quad = r_{bk}.$$

## MAXIMS

Frictions are of great importance in financial markets; they are in many ways the krill that feed the financial Leviathans.

The first step in either finance or mechanics is to consider models that are free of frictions.

Let  $s$  be the sentiment function of gaining wealth. An obvious question is to compare  $s(x)$  with  $s(-x)$ . In other words, how is gaining a hundred pounds different from losing a hundred pounds? Of course, one is making you happier and the other is making you sad but can the happiness and sadness be quantified? and if so, are they equal?

One can argue that  $s$  is linear close to the origin. But what happens as  $|x|$  increases? If you win a billion you are very happy, but if you loose a billion you may very well be broke. You can always gain unlimited amount of money but you always have a limited amount of money to loose, so

$$(14.14) \quad s(x) \neq -s(-x).$$

Another way of thinking about this is the following game. Suppose your annual salary is a hundred pounds and you are invited to play a fair game in which you might win or lose a hundred pounds. Will you play that game? How about a game in which you might win or lose only one pound?

## 15. GLOSSARY

- Systematic Risk: The risk related to the market as a whole. This risk can not be diversified away. Another way to think of this risk is in terms of correlation. The market components are strongly correlated when the market is going through a bad time.
- Non-systematic Risk: Risk unique to an asset.
- Idiosyncratic Risk: Same as Non-systematic risk, unique to the asset.
- CAPM: Capital Asset Pricing Model. The main argument of the model is that the return should depend only on systematic risk. There is a simple equation describing the CAPM model:

$$(15.1) \quad E(r_a) = r_f + \beta(r_m - r_f),$$

where  $r_a$  is the return of the asset,  $r_m$  is the return of the market,  $r_f$  is the risk free rate and  $E(\cdot)$  is the expectation operator.

- Local Vol Model:

## REFERENCES

- [1] Statistics 110 final review. [https://projects.iq.harvard.edu/files/stat110/files/final\\_review.pdf](https://projects.iq.harvard.edu/files/stat110/files/final_review.pdf). Accessed: 2019-10-09.
- [2] Lectures on Smile lecture 01. <http://www.emanuelderman.com/media/smile-lecture1.pdf>. Accessed: 2017-09-05.
- [3] A. B. Dor, L. Dynkin, J. Hyman, P. Houweling, E. van Leeuwen, and O. Penninga. Dtssm (duration times spread). 2007.
- [4] Oxford MAT 2021. <https://www.youtube.com/watch?v=hWxiTz9FdAY>
- [5] A. Papoulis and S. U. Pillai. *Probability, Random Variables, and Stochastic Processes*. Tata McGraw-Hill Education, 2002.
- [6] V. Piterbarg. Funding beyond discounting: collateral agreements and derivatives pricing. *Risk*, 23(2):97, 2010.
- [7] J. M. Steele. *Stochastic Calculus and Financial Applications*. Springer, 2001.

(London) PRET A MANGER, MARBLE ARCH STATION,  
OXFORD STREET, LONDON

Email address, Mohsin Javed: [mhsnjvd@gmail.com](mailto:mhsnjvd@gmail.com)

URL: <http://www.zindajaved.com>



(A) This is the first subcaption



(B) This is the second subcaption

FIGURE 1. This is the overall caption