**Day 1 Summary: Python Core Concepts (Weeks 1-4 Materials)**

- You covered a significant amount of foundational material today, crucial for everything that follows:
- **Core Syntax, Operators & Data Types:**
  - You reviewed arithmetic (+,-,*,/,//,%,**), comparison (>,<,>=,<=,==,!=), and logical (and,or,not) operators.
  - You understand basic data types (int, float, str, bool) and how to check them (type()).
  - You practiced type conversion (int(), float(), str()).
  - Variable assignment (=) and naming rules were covered.
  - **Highlight:** Remember that 10 == 10.0 evaluates to True (int and float equality). Double-check **variable naming rules** (cannot start with a number, no spaces, not a reserved keyword; starting with _ is allowed).
- **Strings:**
  - Creation, special characters (\n, \t), operators (+, *).
  - Indexing (positive/negative) and slicing ([start:stop:step]) including negative steps and reversing.
  - Common methods: len(), .upper(), .lower(), .strip(), .split(), .replace().
  - **Highlight:** Refine the explanation of **string immutability** (the object's content can't change in-place, though the variable can be reassigned). Remember that **string comparison (>, <) is case-sensitive** ("Z" comes before "a"). Ensure you show the output of methods like .replace() when asked. Practice **advanced slicing with negative indices and steps**.
- **Control Flow:**
  - Conditional logic using if/elif/else.
  - Definite iteration using for loops (especially with range()).
  - Indefinite iteration using while loops (requiring initialization and updates).
  - Loop control statements break and continue.
  - The accumulator pattern for summing or counting within loops.
  - **Highlight:** Python's strict **indentation rules** define code blocks – ensure consistency. When writing code, double-check you include all necessary steps (like printing the success message in the password loop).
- **Functions:**
  - Defining functions (def), calling them, understanding the difference between **parameters** (in definition) and **arguments** (in call).

- Using return to send values back; None is returned by default if no return is hit.
- Variable **scope** (local vs. global) and how function calls create local scopes. Assignment inside a function creates a local variable if the name exists globally (unless global keyword is used).
- Docstrings ("""Docstring""") for explanation.
- **Highlight:** Be precise with **function call syntax** (correct name, matching parentheses). Pay close attention to **scope** when predicting output, especially which variable (local or global) is being accessed or printed. Remember intermediate print statements inside functions when predicting output.

- **Lists:**
  - Creation ([]), ordered, mutable (changeable in place).
  - Indexing and slicing (same as strings).
  - Methods like .append(), .sort() (mutates), len(), sum(), min(), max().
  - Looping with for item in my_list.
  - **Highlight:** Key difference between **aliasing (new = old) vs. cloning (new = old[:])** and the side effects. Crucially, review list **operators**: + concatenates elements, * repeats elements; they **do not** create nested lists of the originals. Re-practice **list methods**: pop(index) returns the removed item and modifies the list; remove(value) removes the first matching value and modifies the list; del my_list[index] removes by index. Understand append() vs extend().

- **Dictionaries:**
  - Creation ({}), storing key: value pairs.
  - Keys must be unique and immutable; values can be anything.
  - Accessing values (my_dict[key] or safer my_dict.get(key)). Adding/updating (my_dict[key] = value).
  - Looping through .keys(), .values(), .items().
  - **Highlight:** Practice the **dictionary creation syntax** using colons (:), not equals signs (=). Remember the correct way to **delete items** is typically del my_dict[key]. Update your understanding of order: **Standard Python dictionaries (3.7+) *do* maintain insertion order**, they are not sorted by key automatically.

- **Error Handling:**
  - Using try/except blocks to catch runtime errors.
  - Catching specific errors (ValueError, TypeError, IndexError, KeyError, ZeroDivisionError) vs. bare except:.

- Using else (runs if try succeeds) and finally (always runs).
- **Highlight:** Slightly refine the conceptual understanding of *why* you might raise an exception explicitly (to signal an unrecoverable state immediately).

- **File I/O:**
  - Using with open(…) as f: for safe file handling (auto-close).
  - Modes ('r', 'w', 'a').
  - Reading methods (.read(), .readlines(), iterating for line in f:).
  - Writing methods (.write(), .writelines()).
  - Processing lines: .strip()/.rstrip('\n') for cleaning, .split(delimiter) for parsing.
  - **Highlight:** Master the **looping pattern for reading and processing files line-by-line**. Ensure you correctly handle reading (for line in f:), cleaning (line.strip()), and splitting (line.split(',')) within the loop *before* trying to use the parts. Understand relative vs. absolute paths conceptually.

- **Importing Modules:**
  - Syntax (import x, import x as y, from x import z).
  - Accessing functions/variables based on import type.
  - **Highlight:** Be very careful with the **exact syntax** for imports (from math import pi, pow) and function calls (correct names like mean, correct parentheses). Understand that **namespace pollution** is the primary reason to avoid import *.

**Day 2 Summary: Pandas Power-Up (Weeks 1-4 Materials - Primarily Bootcamp_4_Part_1 & Bootcamp_4_Part_2)**

Today's focus was mastering the core Pandas tools for data manipulation, a critical skill for the exam's coding section.

**Key Concepts Covered:**

1. **Pandas Fundamentals:** Introduced Pandas as the essential library for tabular data, built upon NumPy. Understood the two primary structures: 1D Series and 2D DataFrame.
2. **Series:**
   - Creation from lists and dictionaries.
   - Understanding the index (automatic 0-based default, explicit labels, from dict keys) and accessing data via .values (NumPy array).
   - Indexing/Slicing: Label-based ([], .loc[]), position-based (.iloc[]). Recalled label slicing includes the endpoint.
   - Operations: Applying NumPy functions, Series methods (.mean(), .describe()), arithmetic (creates *new* Series).
   - Filtering: Boolean masking (s[s > value]) and fancy indexing (s[['label1', 'label2']]).
3. **DataFrame Basics:**
   - Concept: Table-like structure, collection of Series sharing an index.
   - Creation from dictionaries of lists or Series (data aligns by index).
   - Attributes: .index, .columns, .values, .shape, .size, .T.
   - Column Selection: df['col'] (Series), df[['col']] (DataFrame), df[['col1', 'col2']] (DataFrame).
4. **DataFrame Modification:**
   - Adding Columns: By assignment (df['new_col'] = ...), including calculations from existing columns.
   - Renaming: .rename(columns={'old':'new'}, inplace=...).
   - Dropping: .drop(labels, axis=0/1, inplace=...). axis=0 for rows (default), axis=1 for columns.
   - Index Manipulation: .set_index('col', inplace=...), .reset_index(drop=False/True, inplace=...).
5. **Indexing & Selection (.loc / .iloc):**
   - Crucial distinction: .loc[] is **label**-based (slice includes endpoint); .iloc[] is **integer position**-based (slice excludes endpoint).

- Syntax: df.loc[row_label(s), col_label(s)], df.iloc[row_pos(es), col_pos(es)]. Can use single labels/positions, lists, or slices.
- Assignment: Can assign values using these indexers (df.loc[...] = value).

6. **Boolean Masking:**
- Creating boolean Series based on conditions (df['col'] > value, df['col'] == 'text').
- Combining conditions with & (and), | (or), ~ (not), using parentheses (cond1) & (cond2).
- Applying masks: df[mask] or df.loc[mask]. Generally creates a **copy**.

7. **Basic Stats & Examination:**
- Aggregations: .describe(), .mean(), .median(), .sum(), .min(), .max(), .count() (non-nulls), .size() (total elements/rows for Series/DF).
- Uniqueness (on Series): .unique(), .nunique().
- Sorting/Ranking: .sort_values(by=...), .nlargest(n, col).
- Metadata: .info() (dtypes, non-null count, memory), .shape.

8. **Grouping (.groupby()):**
- Concept: Split-Apply-Combine.
- Syntax: df.groupby('col_or_list').
- Aggregation: Apply functions like .mean(), .sum(), .min(), .max(), .size() (counts rows per group), .count() (non-nulls per group), .describe() *after* grouping.
- Multi-column/Aggregation: df.groupby(...)['col'].agg_func(), df.groupby(...)[['col1','col2']].agg_func(), df.groupby(...).agg({...}).
- Index: Grouping keys become the index (potentially MultiIndex). .idxmax() finds index label of max value.

9. **Combining DataFrames:**
- pd.concat([df1, df2], axis=0/1, ignore_index=...): Stacking rows (axis=0) or columns (axis=1). Aligns on the *other* axis (outer join by default). ignore_index creates new 0-based index.
- pd.merge(df_left, df_right, on=..., how=..., left_on=..., right_on=..., left_index=..., right_index=...): Database-style joins based on common columns/indices. how determines which keys to keep ('inner', 'outer', 'left', 'right').

10. **File I/O & Plotting:**
- Reading: pd.read_csv(..., delimiter=..., usecols=...), pd.read_excel(..., sheet_name=...).
- Writing: df.to_csv(..., index=False).

- Plotting: Basic .plot(kind=…) method on DataFrames/Series. Seaborn (sns) for enhanced plots (hue argument maps color). Often requires data aggregation *before* plotting grouped results.

**Highlight Areas for Improvement (Based on Quiz Performance):**

- **Series Data Types:** While often homogeneous, remember Series *can* hold mixed types (resulting in object dtype). *(Quiz 1, Q8)*
- **.values vs .array:** The standard attribute for the NumPy array backing a Series/DataFrame is .values. *(Quiz 1, Q4)*
- **.shape Output:** Remember the tuple format is (rows, columns). *(Quiz 2, Q3)*
- **inplace Parameter & Default Behavior:** Methods like .drop(), .rename(), .sort_values() return *new* objects by default. You *must* use inplace=True or reassign the variable (df = df.method(…)) to modify the original DataFrame variable. *(Quiz 3, Q3, Q8)*
- **.rename() Syntax:** Ensure correct dictionary syntax {'old': 'new'}. *(Quiz 3, Q7)*
- **.drop() with Labels:** Use a *list* of labels to drop, not slicing notation, if dropping by label. *(Quiz 3, Q9)*
- **Filtering Application:** Remember the two steps: 1) Create the boolean mask, 2) Apply the mask using df[mask] or df.loc[mask]. *(Quiz 4, Q2 - Initial answer)*
- **View vs. Copy (Filtering/Assignment):** Boolean masking generally returns a *copy*. Modifying this copy won't affect the original. Chained assignment (df[][]=) is discouraged. Use .loc[mask, col] = … for modification based on condition. *(Quiz 4, Q10; Quiz 8, Q5)*
- **Row Count:** Use len(df) or df.shape[0] for the number of rows, not df.count(). *(Quiz 5, Q1)*
- **.unique()/.nunique():** These are *Series* methods, apply them to a selected column df['col'].nunique(). *(Quiz 5, Q4)*
- **.groupby() - .count() vs .size():** Understand .size() counts all rows per group, .count() counts non-null values per column per group. *(Quiz 6, Q3, Q9)*
- **.groupby() - Aggregating Ranges:** Operate on the specific column *after* grouping when calculating ranges or other stats. *(Quiz 6, Q7)*
- **.groupby().describe() Access:** Remember .describe() on grouped data creates a MultiIndex on the *columns*. Access specific stats using tuple indexing: .loc[row_label, (outer_col, inner_stat)]. *(Quiz 6, Q8)*

- **pd.concat() Axis & ignore_index:** axis=0 stacks rows (default), axis=1 stacks columns. ignore_index=True creates a new 0-based index. *(Quiz 7, Q10; Quiz 8, Q10)*. Need to be precise about index result with axis=1 (union of indices). *(Quiz 8, Q2)*
- **pd.merge() - Key Names & Index:** Use left_on/right_on for different key names. Use left_index=True/right_index=True for index merges. *(Quiz 7, Q7, Q8)*
- **pd.read_csv() Arguments:** The argument is usecols, not col. *(Quiz 7, Q5)*
- **.iloc Syntax:** Use lists [] for specific positions, standard slices : for ranges. *(Quiz 8, Q8)*
- **.loc Syntax:** Use square brackets [], not parentheses (). Use actual label type (int vs str). *(Quiz 8, Q9; Quiz 4, Q5)*
- **Plotting Interpretation:** Plotted values reflect the *aggregated* data used (sum, mean, etc.). *(Quiz 7, Q7)*

**Day 3 Summary: Data Concepts, Statistics, and Web Interaction**

Today's focus was broad, moving from fundamental data concepts to how Python interacts with web data sources.

1. **Foundational Concepts:**
   - **Scientific Computing:** Introduced the crucial concepts of **Reproducibility** (getting the same computational result with original data/code) and **Replication** (achieving the same scientific finding, possibly with new data/methods).
   - **Data Types:** Differentiated between structural types (**Cross-Sectional**, **Time Series**, **Panel/Longitudinal**) based on how they handle subjects and time. Also covered the nature of data: **Quantitative** (numerical, measurable/countable, e.g., prices, volume) vs. **Qualitative** (descriptive categories, e.g., text, names, sectors), including **Nominal** (no order) and **Ordinal** (ordered) subtypes.
2. **Financial Calculations & Statistics:**
   - **Returns:** Covered **Simple Returns** (($P_t$ - $P_{t-1}$) / $P_{t-1}$) and **Log Returns** ($\ln(P_t / P_{t-1})$). Understood the benefit of returns (normalization) and the key advantage of log returns (**additivity** over time: total log return = sum of periodic log returns).
   - **Descriptive Statistics:** Reviewed measures of central tendency (**Mean**, **Median**, **Mode**) and dispersion (**Variance**, **Standard Deviation/Volatility**). Briefly touched upon **Correlation Matrices**.
3. **Web Fundamentals:**
   - **Internet vs. WWW:** Differentiated the Internet (global infrastructure, TCP/IP) from the World Wide Web (information system on the internet, uses HTTP, HTML, URLs). Recalled its origins with **Sir Tim Berners-Lee at CERN**.
   - **Client-Server Model:** Understood the basic roles of clients (e.g., browser, script) making requests and servers holding resources and sending responses.
   - **HTTP Protocol:** Detailed the structure of **HTTP Requests** (Request Line: Method/URL/Version; Headers: Host/User-Agent/Accept/*; Body: Optional data for POST) and* ***HTTP Responses** (Status Line: Version/Code/Reason; Headers: Content-Type/Length/*; Body: Requested content). Distinguished key methods (**GET** vs. **POST**) and understood common **Status Code** categories (2xx Success, 4xx Client Error, 5xx Server Error). Contrasted **HTTP**

(unsecured, plain text) with **HTTPS** (secured via **SSL/TLS** encryption and authentication).

4. **Python Web Interaction Tools:**
   - **requests Library:** Practiced basic usage: import requests, requests.get(), requests.post(data=...), passing headers and params, checking response.status_code, accessing response content (.text, .content, .json()), basic try...except error handling, and using time.sleep() for rate limiting.
   - **HTML Basics:** Reviewed core tags (<a>, <h1>, <p>, table tags - <table>, <tr>, <td>, <th>) and attributes (href, id, class) for structuring web content.
   - **BeautifulSoup Library:** Covered parsing HTML (BeautifulSoup(html, 'html.parser')), navigating/searching the structure using direct access (.tag), finding the first element (.find()) or all elements (.find_all()) by tag name and attributes (especially id and class_='...'), extracting text (.text or .get_text()) and attributes (['attr'] or .get('attr')), and the pattern for navigating tables.
   - **API Concepts:** Defined APIs (Application Programming Interfaces) as structured ways for programs to interact. Contrasted their reliability and structured data (JSON/XML) return types with the often less stable HTML scraping approach. Understood the role of API keys and the concept of rate limiting. Recognized that wrapper libraries (like alpha_vantage) simplify API calls.

**Areas Highlighted from Your Quiz/Assignment Performance:**

Based on your answers, these are the specific points to reinforce:

1. **Reproducibility vs. Validity (Quiz 1, Q3 / Assignment Q1):** While you understand the definitions, remember that **reproducibility only verifies the computation, not the underlying scientific truth or generalizability.** A flawed experiment with bad data can still be perfectly reproducible if the code runs correctly on that flawed data.
2. **HTTP vs. TCP/IP Role (Quiz 2, Q3):** Clearly distinguish that **HTTP** is the *application-layer protocol* for the Web (requesting pages/resources), while **TCP/IP** is the underlying *transport/network layer protocol* for the Internet (moving data packets reliably). HTTPS uses HTTP *over* SSL/TLS, which itself runs *over* TCP/IP.
3. **requests Library Specifics:**
   - Accessing response text: Use the attribute response.text, not the method call response.text(). (Quiz 3, Q11).

- Parsing JSON: Use the method response.json(), not json.loads(response). (Quiz 5, Q9).
- POST data argument: Use data=payload for form data (or json=payload for JSON data), not parms=payload. (Assignment Q7, Q8).
- API Key Passing: Keys are usually passed in headers or params according to API docs, not a generic key= argument in requests.get(). (Assignment Q8).
- Error Handling: Need more specific try...except blocks for different requests errors (e.g., requests.exceptions.Timeout, requests.exceptions.ConnectionError) and especially for JSON decoding errors (requests.exceptions.JSONDecodeError or json.JSONDecodeError). (Assignment Q8).
- Status Code Check: if response.status_code: is not standard; use comparisons like if response.status_code == 200: or if 200 <= response.status_code < 300: or the shortcut if response:. (Assignment Q8).

4. **HTTP Concepts:**
   - Host Header: Its purpose is to specify the **domain name** for virtual hosting, not the HTTP version. (Assignment Q1).
   - DELETE Success Code: A successful DELETE usually returns 200 OK or 204 No Content, not 404 Not Found. (Assignment Q5).
   - Cookies: Their primary role is **state management** (sessions, logins, preferences) across stateless HTTP requests. (Assignment Q6).
   - User-Agent: Important for **politeness**, **content negotiation** (mobile vs desktop), and avoiding blocks, not just predefining headers. (Quiz 5, Q3).

5. **BeautifulSoup Specifics:**
   - find() vs find_all() on failure: find() returns None, find_all() returns []. (Assignment Q2).
   - Searching by class: Use the keyword argument class_ (with underscore). (Assignment Q6, Q7).
   - Nested Searching Logic: Need to master finding a container, then finding items *within* that container, often in a loop. (Assignment Q7).
   - Handling Missing Elements: When looping, need to check if find() returned None before trying to access .text or attributes. (Assignment Q7).
   - Tag vs Content: Remember to use .text or .get_text() when you need the content inside the tag, not the tag object itself. (Quiz 5, Q10).

6. **API vs. Scraping Data Format (Quiz 5, Q2):** Reiterate that APIs usually provide structured data (JSON/XML), which is why they are preferred over parsing less structured/stable HTML.

**Day 4 Summary: Market Microstructure & Introduction to Risk/Return**

This revision day covers two main areas: the "rules of the game" in financial markets (Market Microstructure) and the fundamental concepts and calculations for measuring risk and return.

**Part 1: Market Microstructure (Conceptual)**

- **Core Idea:** Understanding how trading actually happens – the mechanisms, participants, costs, and structures.
- **Key Concepts:**
    - **Markets:** Primary (new issues) vs. Secondary (trading existing assets, providing liquidity).
    - **Players:** Exchanges (venues), Brokers/Dealers (intermediaries), Buy-side (investors), Sell-side (service providers).
    - **Liquidity:** Crucial market quality measured by Depth (quantity), Tightness (bid-ask spread), and Resilience (recovery speed).
    - **Trading Mechanisms:** Quote-Driven (dealer prices) vs. Order-Driven (order book matching).
    - **Order Types:** Market (immediate execution, price uncertainty) vs. Limit (price certainty, execution uncertainty).
    - **Electronic Trading:** Evolution from manual; includes Algorithmic Trading (rule-based execution), Direct Access (DMA, Sponsored, Naked), and specialized venues (ATS, Dark Pools, Crossing Networks).
- **Areas Highlighted from Your Quiz Errors:**
    - **Direct Access Nuances (DMA vs. SA):** Need precision on *whose infrastructure* is primarily used for routing (Broker's for DMA, Client's via Broker ID for SA) and where **pre-trade risk controls** primarily reside (Broker for DMA, primarily Client for SA, impacting broker's counterparty risk). *(Q1 7/10, Q1 10/10)*
    - **Algorithmic Trading Limitations/Risks:** Understand the specific weakness of Gen 1 algos (being *static* and predictable, not just input-limited) *(Q2 7/10)*. Recognize the risk of advanced liquidity-seeking algos isn't just execution cost, but **information leakage ("signaling risk")** leading to potential predatory trading. *(Q2 10/10)*
    - **Venue Selection:** For anonymous, large-impact trades, specifically recall **Crossing Systems/Dark Pools** as key mechanisms mentioned, rather than more general terms like OTC. *(Q4 7/10)*

- **Naked Access Risk:** The primary risk is the broker's **lack of pre-trade control/filtering** over orders hitting the market under their name. *(Q8 7/10)*

## Part 2: Introduction to Risk & Return (Conceptual & Basic Python)

- **Core Idea:** Quantifying potential rewards (returns) and uncertainty (risk) associated with financial assets.
- **Key Concepts:**
    - **Returns:** Simple (pct_change()) vs. Log (time-additive, preferred for multi-period). Expected Return (often proxied by historical mean).
    - **Risk (Volatility):** Variance (dispersion, squared units) vs. Standard Deviation (dispersion, same units as return, more interpretable for single assets).
    - **Risk (Co-movement):** Covariance (measures joint variability, scale dependent) vs. Correlation (standardized -1 to +1, measures strength and direction of *linear* relationship). Crucial for diversification.
    - **Annualization:** Scaling periodic measures (daily) to annual. *Key difference:* Multiply mean log return & variance by periods (e.g., 252); multiply standard deviation by sqrt(periods).
    - **Diversification:** Combining assets with low/negative correlation reduces portfolio risk (unsystematic component) because assets don't move perfectly together.
    - **Portfolio Risk:** Depends *critically* on how assets co-move (covariance/correlation), not just a weighted average of individual variances.
- **Areas Highlighted from Your Quiz Errors:**
    - **Covariance vs. Correlation Definition:** Must be precise. Correlation is **standardized** (-1 to +1) measuring **strength** and direction; Covariance scale is hard to interpret directly. *(Q3 7/10)*
    - **Input Data:** Remember the minimum required input is a **time series of prices**. *(Q5 7/10 - Skipped)*
    - **Portfolio Variance Calculation:** Understand *why* covariance is included – it captures the **interaction effect** between assets. *(Q6 10/10)*

## Part 3: Python Implementation

- **Core Idea:** Using standard libraries (Pandas, NumPy) to perform risk/return calculations.
- **Key Functions:**
    - Returns: .pct_change(), np.log(), .shift(), .diff()

- Descriptive Stats: .mean(), .var(), .std()
- Co-movement: .cov(), .corr()
- Other: np.sqrt(), .rolling()
- **Areas Highlighted from Your Quiz Errors:**
  - **Log Return Formula:** Correct calculation is np.log(prices_df / prices_df.shift(1)) or np.log(prices_df).diff(). Subtraction before the log is incorrect. *(Q7 7/10)*
  - **Portfolio Variance/Std Dev Code:** The correct approach uses matrix multiplication: variance = weights.T @ annual_cov_matrix @ weights, followed by np.sqrt() for std dev. Simple weighted sums of std dev/covariance are incorrect. *(Q8 10/10)*
  - **Extracting Specific Cov/Corr Values:** Need to use .loc[row, col] or specific methods like series1.cov(series2) to get values between specific assets from the full matrix generated by .cov() or .corr(). *(Q9 7/10)*
  - **Annualizing Rolling Calculations:** Remember to apply the annualization factor (e.g., * np.sqrt(252)) *after* the rolling calculation (e.g., .rolling(window).std()). *(Q10 10/10)*

**Day 5 Summary: Portfolio Theory & Application**

Today's focus was on moving from individual asset risk/return to understanding how combining assets into portfolios works, guided by Modern Portfolio Theory (MPT). The goal is to construct "optimal" portfolios that balance risk and return.

**Key Topics Covered:**

1. **Portfolio Fundamentals:**
   - **Definition:** A portfolio is a collection of assets.
   - **Weights:** Represent the proportion of total investment value in each asset. Weights sum to 1 for standard long-only portfolios, can sum to 0 (dollar-neutral), and individual weights can be negative (shorting) or >1 (leverage).
   - **Expected Return:** The portfolio's expected return is the weighted average of the individual assets' expected returns (E[Rp] = w.T @ E[R]).

2. **Portfolio Risk (Variance/Volatility):**
   - **Crucial Concept:** Portfolio risk (variance/volatility) is **NOT** the simple weighted average of individual asset risks.
   - **Covariance/Correlation:** It depends critically on how assets move together, measured by their **covariance** or **correlation**. Lower correlation between assets leads to greater risk reduction benefits.
   - **Calculation:** Requires accounting for all individual variances *and* all pairwise covariances. The matrix formula is key: $\sigma p^2$ = w.T @ Σ @ w.
   - **Diversification:** Combining imperfectly correlated assets reduces *unsystematic* (firm-specific) risk.

3. **Python Implementation:**
   - Calculating daily log returns from prices (np.log(prices / prices.shift(1))).
   - Annualizing mean returns (daily_mean * 252) and the covariance matrix (daily_cov * 252).
   - Calculating portfolio expected return (np.dot(weights, annual_means)).
   - Calculating portfolio variance (weights.T @ annual_cov_matrix @ weights or using np.dot).
   - Calculating portfolio volatility (np.sqrt(portfolio_variance)).

4. **Modern Portfolio Theory (MPT) Concepts:**
   - **Objective:** Maximize expected return for a given level of risk, or minimize risk for a given level of expected return.
   - **Efficient Frontier:** The set of portfolios offering the best possible expected return for each level of risk. Portfolios below the frontier are suboptimal.

- **Risk-Free Asset:** Introduces a straight line (Capital Market Line) representing combinations of the risk-free asset and the optimal risky portfolio.
- **Tangency Portfolio:** The specific portfolio of *risky assets* on the original efficient frontier that offers the highest Sharpe Ratio when combined with the risk-free asset. MPT suggests all investors should hold this risky mix.
- **Sharpe Ratio:** Measures excess return (above risk-free rate) per unit of total volatility ((E[Rp] - Rf) / σp). Higher is better. The Tangency Portfolio maximizes this ratio.

5. **Advanced Concepts:**
   - Impact of constraints (like no short selling) on the feasible region and efficient frontier.
   - Limitations of MPT, especially reliance on historical estimates (non-stationarity, non-normality/fat tails).
   - Specific conditions for achieving zero risk with two assets.
   - Formulating optimization problems (e.g., minimizing negative Sharpe Ratio).

**Highlights from Your Quiz Answers (Areas for Review):**

- **Weight Interpretation (Fundamentals Quiz Q1):** While you mentioned dollar-neutral, ensure you can also clearly state that *individual* negative weights mean shorting, and weights >1 (or sum >100%) imply leverage.
- **Diversification Limit Explanation (Risk Quiz Q3):** The primary reason portfolio risk doesn't reach zero is **systematic risk** (market risk), which affects all assets and cannot be diversified away. Your answer mentioned other risks (liquidity), which wasn't the core MPT reason.
- **Calculation Accuracy (Risk Quiz Q4a):** A minor slip in substituting the correct volatility value (used 30% instead of 35%) into the variance formula. Double-check calculations.
- **Zero-Risk Portfolio Weights (Advanced Quiz Q3):** Achieving zero risk requires ρ = -1, but the specific weights depend on the *relative volatilities* ($w_A = \sigma_B / (\sigma_A + \sigma_B)$, $w_B = \sigma_A / (\sigma_A + \sigma_B)$), not necessarily 50/50.
- **Optimization Formulation (Advanced Quiz Q4):** You need to be able to state the function to *minimize* (e.g., -Sharpe Ratio) and the standard constraints (weights sum to 1, potentially non-negative weights) when setting up an optimization problem.
- **CRITICAL Python Variance Calculation (Advanced Quiz Q5a): This was the most significant error.** You used mean_returns_annual inside the portfolio variance

calculation (np.dot(weights.T, np.dot(MEAN_RETURNS, weights))). The formula **must** use the **cov_matrix_annual**: np.sqrt(np.dot(weights.T, np.dot(COVARIANCE_MATRIX, weights))) to get volatility. Ensure you firmly distinguish between inputs for expected return vs. variance/volatility calculations.

**Summary of Day 6 Revision:**

Day 6 covered a broad range of conceptual and practical topics crucial for understanding modern data analysis techniques in finance and accounting:

1. **Machine Learning Concepts:**
   - Defined the hierarchy of AI, ML, and Deep Learning.
   - Distinguished between Structured and Unstructured data.
   - Explained Supervised (Classification, Regression) vs. Unsupervised (Clustering, Dimensionality Reduction) learning paradigms.
   - Defined "Features" as input variables for models.
   - Explained "Overfitting" and its problems.
   - Introduced the Confusion Matrix (TP, FP, FN, TN) for classification evaluation.
   - Defined key performance metrics: Sensitivity (Recall), Specificity, and Precision (PPV).
2. **PCA & Factor Analysis:**
   - Understood their purpose as data reduction techniques.
   - Distinguished between PCA (variance maximization) and Factor Analysis (latent factor identification).
   - Reviewed PCA methodology: role of standardization, correlation matrix, eigenvalues (variance explained), eigenvectors (weights).
   - Learned interpretation: Factor/Component Loadings (correlation between variable and component), naming components.
   - Studied component retention methods: Kaiser's Rule ($\lambda > 1$), Scree Plots (elbow method).
   - Understood Component Rotation (Orthogonal vs. Oblique) to improve interpretability.
   - Reviewed data suitability tests: KMO (sampling adequacy) and Bartlett's Test (presence of correlations).
3. **Data Ethics:**
   - Defined Data Ethics, focusing on personal data.
   - Reviewed principles of Informed Consent.
   - Discussed Data Ownership vs. individual control/privacy rights.
   - Explored Privacy challenges (persistence, anonymity, boundaries, Right to be Forgotten).

- Understood the Collection vs. Use distinction regarding harm.
- Identified common issues (permission creep, digital footprints, de-identification challenges).
- Named Leakage Types (identity, attribute, link, membership).
- Recognized the ethical implications of Data Validity (using bad data/models).

4. **Textual Analysis & Large Language Models (LLMs):**
   - Distinguished between Textual Analysis (counts, presence) and NLP (meaning, structure, context).
   - Acknowledged diverse Text Data sources.
   - Understood Readability concepts and the Gunning Fog Index.
   - Defined LLMs by scale (parameters, data).
   - Explained Parameters conceptually (weights storing knowledge).
   - Reviewed the two-stage Training Process (Pre-training -> Base Model; Fine-tuning -> Assistant Model).
   - Understood the core "predict next word" objective and its role in knowledge acquisition.
   - Briefly touched on limitations (e.g., reversal curse) and future concepts (System 1/2, RAG).

5. **Practical Python Coding (PDF/Readability):**
   - Identified key libraries: PyPDF2, pdfminer.six, textstat.
   - Practiced core tasks: opening PDFs, getting page counts, accessing pages, extracting text (page-by-page with PyPDF2, whole file with pdfminer.six).
   - Practiced calculating Gunning Fog index with textstat.
   - Reviewed looping through pages (PyPDF2) and joining text.

**Highlights of Areas Needing Review (Based on Quiz Performance):**

While you showed a good grasp of many core concepts, the quizzes highlighted specific areas needing more attention, particularly for achieving that 10/10 difficulty level:

1. **Precise Interpretation of ML Metrics:**
   - You need to be very clear on *exactly* what **Sensitivity (Recall)**, **Specificity**, and **Precision (PPV)** measure. (Errors in Q1c, Q1d, Q1e of the final quiz; Q4a, Q4b, Q5a, Q5b, Q5c of the first quiz).
   - Crucially, understand **which metric aligns with specific goals** (e.g., minimizing False Negatives requires high Sensitivity; minimizing False

Positives often relates to high Specificity or Precision depending on the context). (Errors in Q1d, Q1e).

2. **PCA Calculations & Techniques:**
   - Remember the correct denominator for calculating **% variance explained** from eigenvalues (it's the *total number of original variables*, not just the sum of listed eigenvalues). (Errors in Q3a, Q3b of PCA quiz).
   - Clearly distinguish between **pre-analysis suitability tests (KMO, Bartlett's)** and **post-extraction interpretation techniques (Rotation)**. (Errors in Q6a, Q6b of PCA quiz).
   - Know the difference and purpose of **Orthogonal (e.g., Varimax) vs. Oblique (e.g., Promax) rotation**. (Error in Q6b).

3. **Nuanced Ethical Reasoning:**
   - Go beyond surface-level answers. Elaborate on *why* something is an ethical issue (e.g., *why* overfitting is ethically problematic in loan decisions, not just that it's inaccurate). (Weakness in Q1a of final quiz).
   - Consider multiple perspectives and the specific context when evaluating consent and purpose limitation (e.g., LLM summarization Q8).

4. **Data Type Identification:**
   - Be certain about classifying data types, especially recognizing that raw **text is unstructured**. (Error in Q1b of first quiz).

5. **Python Error Handling & Precision:**
   - While you grasp the basic library usage, focus on writing more **robust code**, particularly including specific **error handling** (try-except blocks for file/page errors). (Weakness in Q6 of final quiz).
   - Pay attention to **syntax details** (like f-strings). (Minor error in Q2 of Python quiz).

**Week 7: Market Microstructure (Take Home Lecture Notes) - Extracted Content**

1. **Assets: Stocks, Bonds, Options, Forex, ETF, Cryptocurrency**
   - **Stocks (Slide 5):** (Also known as shares or equity) is a type of security that signifies proportionate ownership in the issuing corporation. It entitles the holders to that proportion of the corporation's assets and earnings.
   - **Bonds (Slide 6):** A fixed income instrument that represents a loan made by an investor to a borrower.
   - **Options (Slide 7):** An agreement between a buyer and seller that gives the purchaser of the option the right to buy or sell a particular asset at a later date at an agreed price. Most often used in securities, commodities, and real estate transactions.
   - **FOREX (Slide 8):** Known as foreign exchange or currency trading is a decentralized global market where all the world's currencies trade. It is the most liquid market in the world with an average daily trading volume exceeding $5 trillion. Active 24/7 due to different time zones.
   - **Cryptocurrencies (Slide 9):** A digital asset designed to work as a medium of exchange. Decentralized assets currently not regulated or controlled by any one country, central bank, or regulatory authority (except El Salvador mentioned as an exception).
   - **ETF (Exchange-Traded Fund) (Slide 19):** (Mentioned under Institutional Trading Types, relevant as an asset class) An example of a basket of securities that trade on an exchange just like a stock. ETFs are a collection of stocks or securities that track an underlying index.
   - *(Note: "Funds" were listed on Slide 4 but not detailed further in the overview slides.)*
2. **Trading Frequency (Slide 54):** The determinant of when requirements are turned into executions.
   - **Continuous trading:** Leads to price volatility.
   - **Periodic Trading:** Scheduled for specific time/s in the day.
   - **Request-driven trading:** Requesting a quote from a market maker, but not the most efficient in terms of the price achieved.
3. **Fill Instructions: Immediate-or-Cancel (IOC), Fill or Kill (FOK), All-or-None (AON)**
   - **Content NOT FOUND:** These specific fill instruction types (IOC, FOK, AON) are **not mentioned** in the provided Week 7 lecture notes text file.
4. **Order Concepts: Preference & Directed Orders**
   - **Content NOT FOUND:** These specific order concepts (Preference Orders, Directed Orders) are **not mentioned** in the provided Week 7 lecture notes text file.
5. **Routing Instructions: Do-not-route, Directed-routing, Intermarket sweep**
   - **Content NOT FOUND:** These specific routing instruction types (Do-not-route, Directed-routing, Intermarket sweep) are **not mentioned** in the provided Week 7 lecture notes text file. (While order routing is discussed conceptually (e.g., Slides 15, 30, 33, 37), these specific instructions are absent).
6. **Transaction Costs: Investment-related & Trading-related**

- **Partial Content Found (Slide 28):** The text discusses "transaction cost analysis" which breaks down various costs. It highlights that factors beyond just market impact are significant:
  - **Timing risk**
  - **Opportunity costs**
  - These can outweigh **market impact**.
- *(Note: The text doesn't provide a distinct list differentiating "Investment-related" (like Taxes, Delay) and "Trading-related" (like Spreads, Price trend) as per the original checklist structure, but it does cover key components of trading-related costs within the context of algorithm evolution.)*

7. **Trade Analysis: Pre-trade vs. Post-trade Analysis**
   - **Pre-trade Analysis (Mentioned):**
     - **Slide 34 (Sponsored Access):** Mentions brokers usually need to monitor trading pre-trade (using fast systems or solutions) to ensure no excessive risks are taken. Contrasts this with "Naked Access" where monitoring is post-trade.
   - **Post-trade Analysis (Implied):**
     - **Slide 28 (Transaction Cost Analysis):** Discusses the application of transaction cost analysis, which inherently involves analyzing costs *after* a trade has occurred to understand factors like market impact, timing risk, and opportunity costs. The term "Post-trade analysis" isn't explicitly used here, but the concept is present.

**Extracted Content for Missing Checklist Items (Weeks 7-9)**

1. **Excess Return (Definition)**
   o **From Week 7 Slides (Slide 3):**
     ▪ Formula: Excess Return = Rit - rf
     ▪ Context: "We're going to be looking at excess returns, which is in excess of the net risk-free rate".
     ▪ Relation to Risk Premium: "So the excess return you can think of as a realization of that risk premium."

2. **Risk Premium (Definition)**
   o **From Week 7 Slides (Slide 3):**
     ▪ Formula: Risk Premium = E[Rit] - rf
     ▪ Context: "The average rate of return of a risky security minus the risk-free rate".
     ▪ Context: "...the number that we're going to be concerned with most is this risk premium number."
   o **Implicitly in Week 7 Notebook (Beta section):** The term [ E (Rm) - Rf ] in the CAPM formula E (Rp) = Rf + Bp [ E (Rm) - Rf ] represents the *Market Risk Premium*.

3. **Market Properties: Efficient Markets (Concept)**
   o **From Week 7 Slides (Slide 7):**
     ▪ Headline: "What Properties Should Stock Prices Have in 'Efficient' Markets?"
     ▪ Properties Listed:
       ▪ Random, unpredictable.
       ▪ Prices should react quickly and correctly to new information.
       ▪ Investors cannot earn abnormal, risk-adjusted returns (in other words, once risk adjustment is taken into account, there shouldn't be any additional return left over.)
     ▪ Context: "That's what we think of as a well-functioning & highly competitive market."

4. **Common Risk Measures: Alpha (α), Beta (β) (Definitions/Formulas)**
   o **Alpha (α):**
     ▪ **From Week 7 Slides (Slide 11):** "the active return on an investment, or an excess return (abnormal rate of return) above the benchmark. or the portion of the excess return that is NOT explained by systematic risk".
     ▪ **From Week 7 Notebook (Beta section):** Represented as the intercept term (α) in the regression equation Rp = α + Bp * Rm + ε. The notebook code calculates this const term.
   o **Beta (β):**
     ▪ **From Week 7 Slides (Slide 11):** "beta measures the amount of systematic risk an individual security has relative to the whole market (e.g., S&P 500)."
     ▪ Formula (Coefficient): Beta coefficient (β) = cov(Rs, Rm) / var(Rm)

- Interpretation: "The market has a beta of 1"; "A security with a beta>1(<1) indicates that it is more(less) volatile than the stock market." (Further interpretation examples on Slide 13).
        - **From Week 7 Notebook (Beta section):** Calculated as the slope coefficient (Bp) in the regression $Rp = \alpha + Bp * Rm + \varepsilon$. The notebook code extracts this coefficient using statsmodels.
5. **Concept: Risk-Return Tradeoff**
    o **From Week 7 Slides (Slide 14):**
        - Definition: "The risk-return tradeoff states that the potential return rises with an increase in risk."
        - Application: "Investors often use the risk-return tradeoff as one of the essential components of each investment decision, and it is also used to assess the portfolio."
    o **(Visualized extensively in Week 8 Slides (e.g., Slide 23, 25, 27, 28) through Mean-Variance plots, but defined explicitly in Week 7.)**
6. **Portfolio Types (Conceptual - Beyond weight definitions)**
    o **From Week 7 Slides (Slide 20):** Specific types defined:
        - **Income portfolio:** Securing a steady flow of income (e.g., dividend stocks) rather than focusing on price appreciation.
        - **Growth portfolio:** Parks money into growth stocks (companies in active growth stage), potentially higher risk.
        - **Value portfolio:** Puts money into assets considered cheap in valuation, focusing on bargains. Often sought during economic recessions.
    o **(Other types mentioned by weight structure in Week 7/8 Slides: Long-only (weights sum to 1), Dollar-neutral (weights sum to 0), Leveraged (weights sum > 1 or individual > 1), Short positions (negative weights))**

**Weeks 11-12: Textual Analysis & LLM Concepts - Extracted Content**

1. **Neural Network (Basic Concept)**
   - **From Week 11 Textual Analysis & LLM - Part 1 (Slide 9):**
     - LLMs like Llama-2 utilize a **Neural Network (NN) architecture**. The weights (parameters) of the model *are* the weights of this neural network. A large number of parameters (e.g., 70 billion for Llama-2 70b, taking ~140GB) constitute the NN.
     - Conceptually, the NN architecture itself can be implemented with relatively simple code (e.g., 500 lines of C mentioned for Llama), but the complexity lies in the vast number of parameters (weights) that define its specific behaviour.
   - **From Week 11 Textual Analysis & LLM - Part 1 (Slide 11):**
     - **Function:** In the context of LLMs, the Neural Network's primary function discussed is **predicting the next word** in a sequence. It takes a sequence of input words (context) and outputs a probability distribution for the next word (e.g., "mat" with 97% probability after "the cat sat on a").
     - **Structure:** Shows a visual diagram of interconnected nodes arranged in layers (input layer, hidden layers, output layer), representing a typical NN structure.
     - **Parameters:** Billions of parameters (weights) are dispersed throughout the entire NN. These parameters are adjusted iteratively during training to improve the network's performance on the next-word prediction task. The exact role of each individual parameter among the billions is generally not fully understood ("we don't actually really know what these billions of parameters are doing").
     - **Training:** The process of training the NN is linked to compression; the ability to predict the next word efficiently implies a form of compressing the information from the vast training data (like the internet) into the network's parameters.
   - **From Week 11 Textual Analysis & LLM - Part 1 (Slide 13):**
     - Shows a schematic diagram specifically for a **transformer NN architecture**, commonly used in modern LLMs.
     - Reiterates that parameters are dispersed through the network and adjusted (optimized) iteratively during training for better next-word prediction.
   - **From Week 11 Textual Analysis & LLM - Part 2 (Slide 2):**
     - Mentions Stage 1 (Pre-training) involves compressing vast amounts of text (~10TB) into a **neural network**.

**Supplementary Content for Missing Checklist Items (Week 7 Scope)**

**(Difficulty Assessment: Low to Medium - These are standard market terminology, slightly more specific than basic Market/Limit orders but fundamental to execution.)**

1. **Fill Instructions: Immediate-or-Cancel (IOC), Fill or Kill (FOK), All-or-None (AON)**
   o These instructions provide additional conditions regarding the time and quantity constraints for an order's execution, typically used with limit orders.
   o **Immediate-or-Cancel (IOC):**
      ▪ Requires that all or part of the order must be executed *immediately* upon being received by the market.
      ▪ Any portion of the order that cannot be filled immediately is cancelled. Partial fills are allowed.
      ▪ *Purpose:* To grab available liquidity instantly without leaving an unfilled order resting on the book.
   o **Fill or Kill (FOK):**
      ▪ Requires that the *entire* quantity of the order must be executed *immediately*.
      ▪ If the entire quantity cannot be filled immediately, the entire order is cancelled (killed). No partial fills are allowed.
      ▪ *Purpose:* Ensures the full desired quantity is executed right away, or not at all; avoids partial fills.
   o **All-or-None (AON):**
      ▪ Requires that the *entire* quantity of the order must be executed, but *not necessarily immediately*.
      ▪ The order can rest in the order book until the full quantity is available at the specified limit price (or better). No partial fills are allowed.
      ▪ *Purpose:* To ensure the full position size is achieved, preventing partial fills, often used for specific strategies or block trades where a partial fill is undesirable. Main difference from FOK is the lack of an immediacy requirement.
2. **Order Concepts: Preference & Directed Orders**
   o These concepts relate to specifying the initial destination or handler for an order.
   o **Preference Orders:** (Primarily a concept in some markets, e.g., historically in the US for retail orders)
      ▪ An instruction from a broker (often on behalf of a retail client) to route an order to a *specific dealer or market maker*, potentially with the expectation of price improvement or other execution quality benefits. The dealer is "preferred."
   o **Directed Orders:**
      ▪ An instruction given by the trader/client to their broker specifying the *particular exchange or trading venue* (e.g., NYSE, LSE, specific ATS) to which the order should be sent first.
      ▪ *Purpose:* Gives the trader control over the initial placement venue, perhaps due to perceived liquidity, fees, or strategic reasons.

3. **Routing Instructions: Do-not-route, Directed-routing, Intermarket sweep**
   - These instructions govern how an order interacts (or doesn't interact) with multiple trading venues in a fragmented market system.
   - **Do-not-route (DNR / DNT):**
     - An instruction attached to an order (usually a limit order) telling the receiving exchange or venue *not* to route the order to any other trading center, even if a better price might be available elsewhere.
     - *Purpose:* Often used when a trader wants to *add* liquidity to a specific venue or believes their order will be filled at the desired price on that venue without needing to search elsewhere. It contrasts with standard "smart order routing" where the broker actively seeks the best price across venues.
   - **Directed-routing:** *(While sometimes used, it largely overlaps with "Directed Orders" concept above or falls under broker's specific routing logic rather than a standard tag like DNR/ISO. The key contrast is DNR vs. standard routing.)*
   - **Intermarket Sweep Order (ISO):**
     - A specialized type of limit order (or multiple limit orders) used to execute rapidly against displayed liquidity across *several different trading venues simultaneously*.
     - It's designed to comply with regulations (like Reg NMS in the US) that generally prohibit "trading through" a better-priced protected quote on another market. By "sweeping" multiple markets at once up to its limit price, an ISO can access liquidity quickly without violating these rules.
     - *Purpose:* Primarily used to rapidly take liquidity from multiple market centers at once, often by sophisticated or algorithmic traders needing immediate execution across the market landscape.