

Introduction to Data Science

Course 094201

Lab 2:

Basic C++ and Data Analysis

Spring 2017

- Copy the following to your local folder and unzip plots.zip :
`cp -r /mnt/share/students/LAB2/lab2_for_students your_folder`
`cp /mnt/share/students/LAB2/plots.zip your_folder`
`unzip your_folder/plots.zip`
- CLion will be our workspace for C++ this semester
- Open CLion (Applications->Programming->CLion)
- Open the project you just copied.
- Compile the project with CLion.
- Open the folder of the project and find the executable file.
- Run the executable outside of CLion (like in lab1).

The Data

- רקע:
- לרשותכם קובץ נתונים על איכות היין.
- הנכם מעוניינים לבצע חקר של הקשרים בין המשתנים השונים בקובץ: כגון רמת חומציות, רמת הסוכר והאלכוהול וכו'. איכות היין (המשתנה quality) היא מדד סובייקטיבי הניתן על סמך שיפוט אנושי.
- שימו לב שכאשר מדובר בבעיות חיזוי סטטיסטיקאים נוטים לחלק את המשתנים (variables) למשתנה תלוי (dependent) - שזה אותו משתנה שמנסים לחזות ומשתנים מסבירים (explanatory). למשל אם המטרה הייתה חיזוי איכות היין, אזי המשתנה של האיכות היה התלוי ומשתנים אחרים בקובץ היו המסבירים. לעיתים נקרא המשתנה התלוי גם response variable.
- אנשי Machine Learning משתמשים במינוח אחר (בד"כ) למשתנה התלוי קוראים class או label או outcome variable. ואילו למשתנים המסבירים קוראים: features (תכונות) או predictors. כאמור לעיתים משתמשים במושגים אלו במעורב.

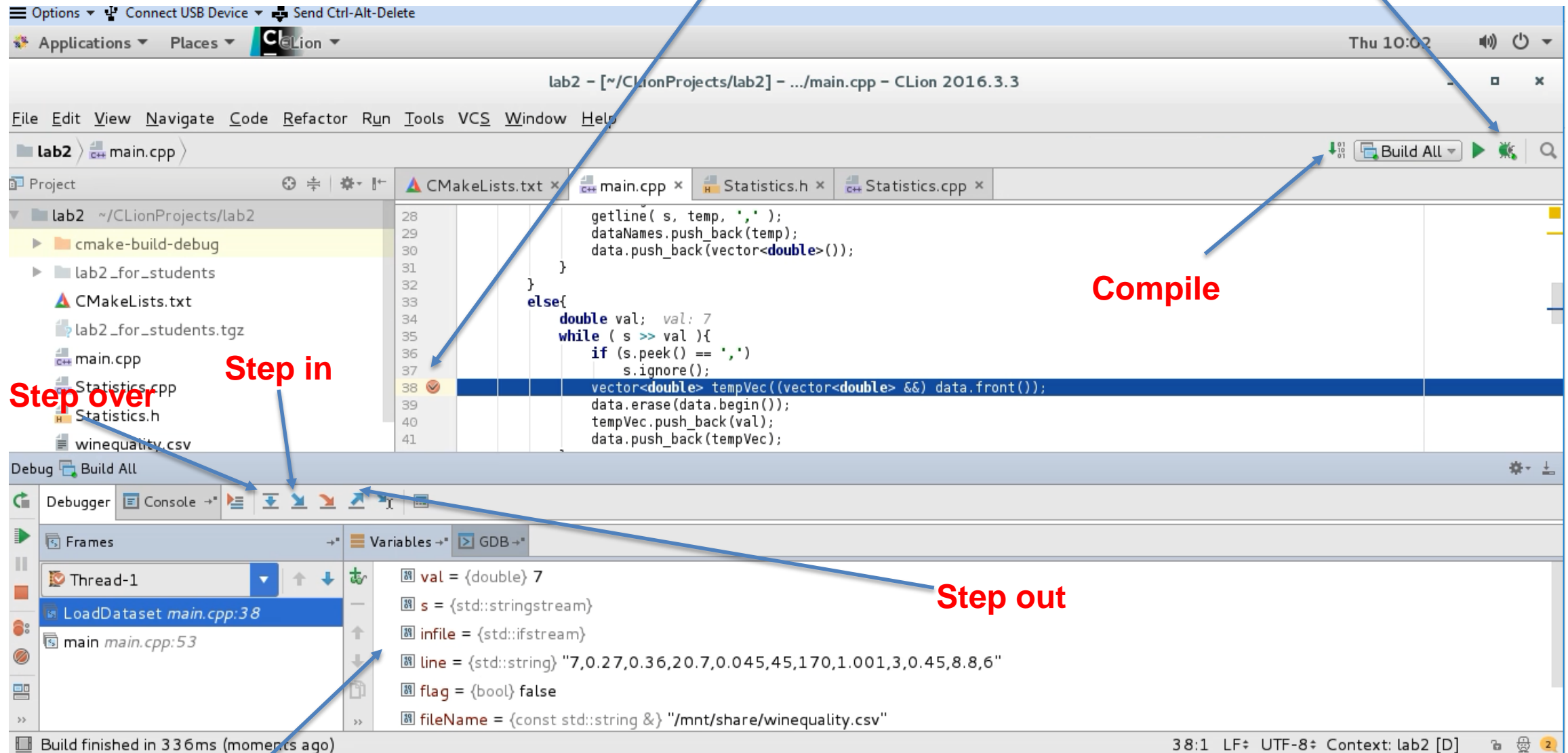
Reading the data

- בטרם אתם עובדים עם קובץ נתונים כלשהו מומלץ לאסוף עליו כמה שיותר מידע
- הסתכלות על השורות הראשונות והאחרונות של הנתונים כבר נותנת לנו קצת מידע על פורמט הקובץ (מידע מצומצם כמובן)
- אתם יכולים לעשות זאת על ידי שימוש בפקודות `bash`: `tail` ו `head` שלמדנו במעבדה הקודמת
- בקובץ `main.cpp` מצויה פונקציה שטוענת את קובץ הנתונים
- ראו שאתם מבינים מהם מבני הנתונים שבהם משתמשים בתרגיל
- שימו לב לאופן העברת הפרמטרים לפונקציה
- הבינו מתי הפונקציה מחזירה `true` ומתי היא מחזירה `false` ולמה הבדיקות הללו נחוצות (ולמה זה לא מספק עבור קובץ חדש)

Debug with CLion

Break point

debug



Compile

Step over

Step in

Step out

Variables

Debugging with CLion

- Let's see what happens if we want to stop the run in middle and see the current state of our program.
- **Set a few “break points” in main.cpp and debug the program**
- **Set a break points to see :**
 1. The value of the 4th feature in the 3th iteration.

 2. The name of the 6th feature.

- The files Statistics.h and .cpp should contain some statistical measures we learned at class
- Implementation of the Median function is given to you, your task is to complete the other measures
- The Median function contains call to standard library (STL) algorithm: sort, we learn more STL algorithms later on

- Do it yourself: implement the following:
 - Mean: a function that receives a const reference to vector and returns its mean.

Mean (\bar{x})

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Do it yourself: implement the following:
 - Variance: a function that receives a const reference to vector and returns its variance.

The **sample variance**:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$