# Pyspark Install

## Introduction to Massive Data Analysis

2021.09

# Step 1 - Python Install and Setting

## Download Python files

- Available download page: https://www.python.org/downloads/

- Recommended version: python3.5.x, python 3.6.x

- It may cause some problems while downloading by Anaconda

# Step 1 - Python Install
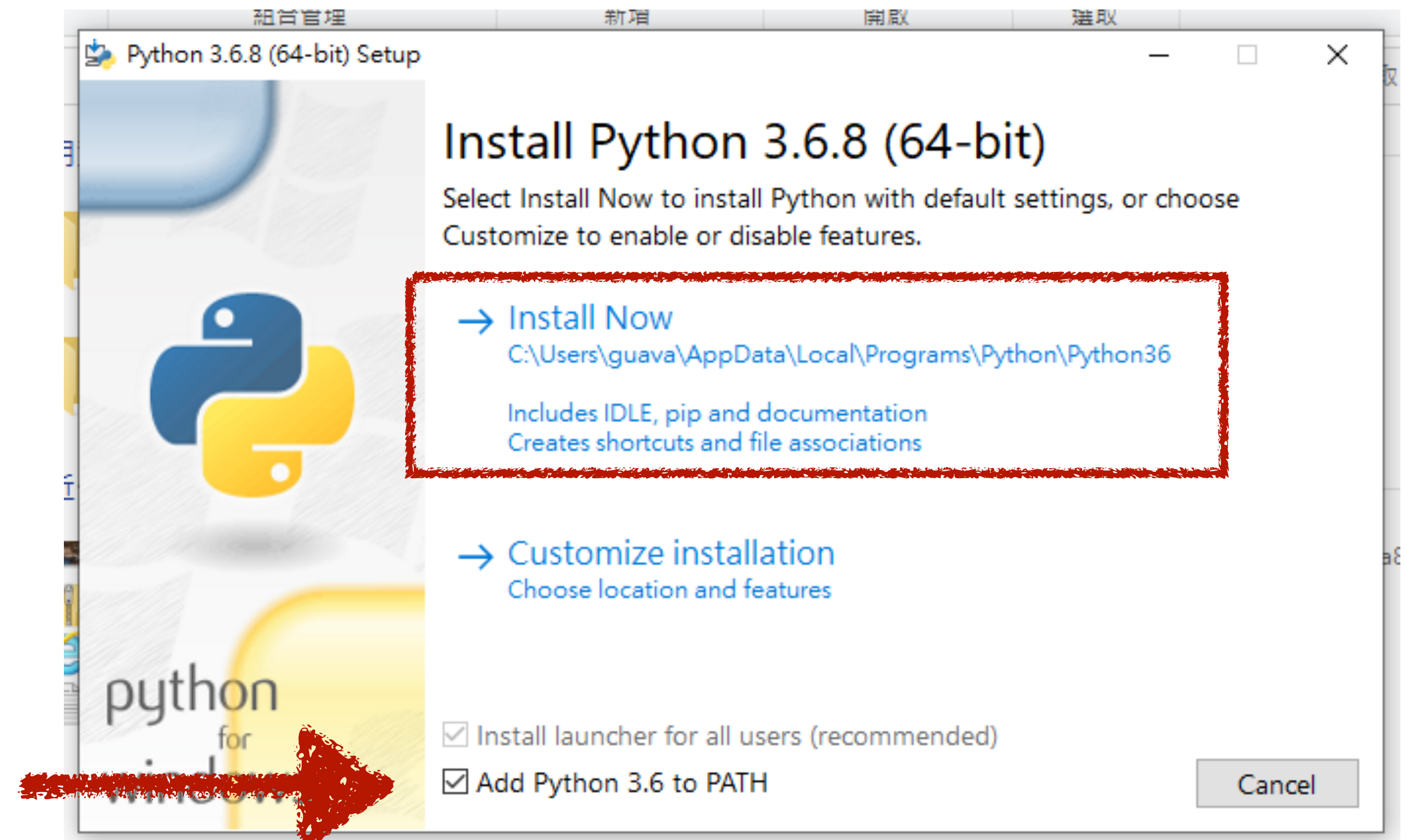## Download Python files

- Select **web-based installer**

# Step 1 - Python Install

## Installing and setting environment variables

- Unzip and run the downloaded file

- **Select "Add to PATH"**

- Install now

# Step 2 - Java Install and Setting
## Download Java jdk installer

– Available download page: https://www.oracle.com/java/technologies/downloads/#java8-windows

– Should select **Java8**

# Step 2 - Java Install and Setting
## Setting environment variables

- At Environment Variables > System Variables

- Add **JAVA_HOME**

- Fill in the **actual address** of your jdk folder
  (not the one in picture)

# Step 2 - Java Install and Setting
## Setting environment variables

- At Environment Variables > System Variables

- Modify and add **PATH**

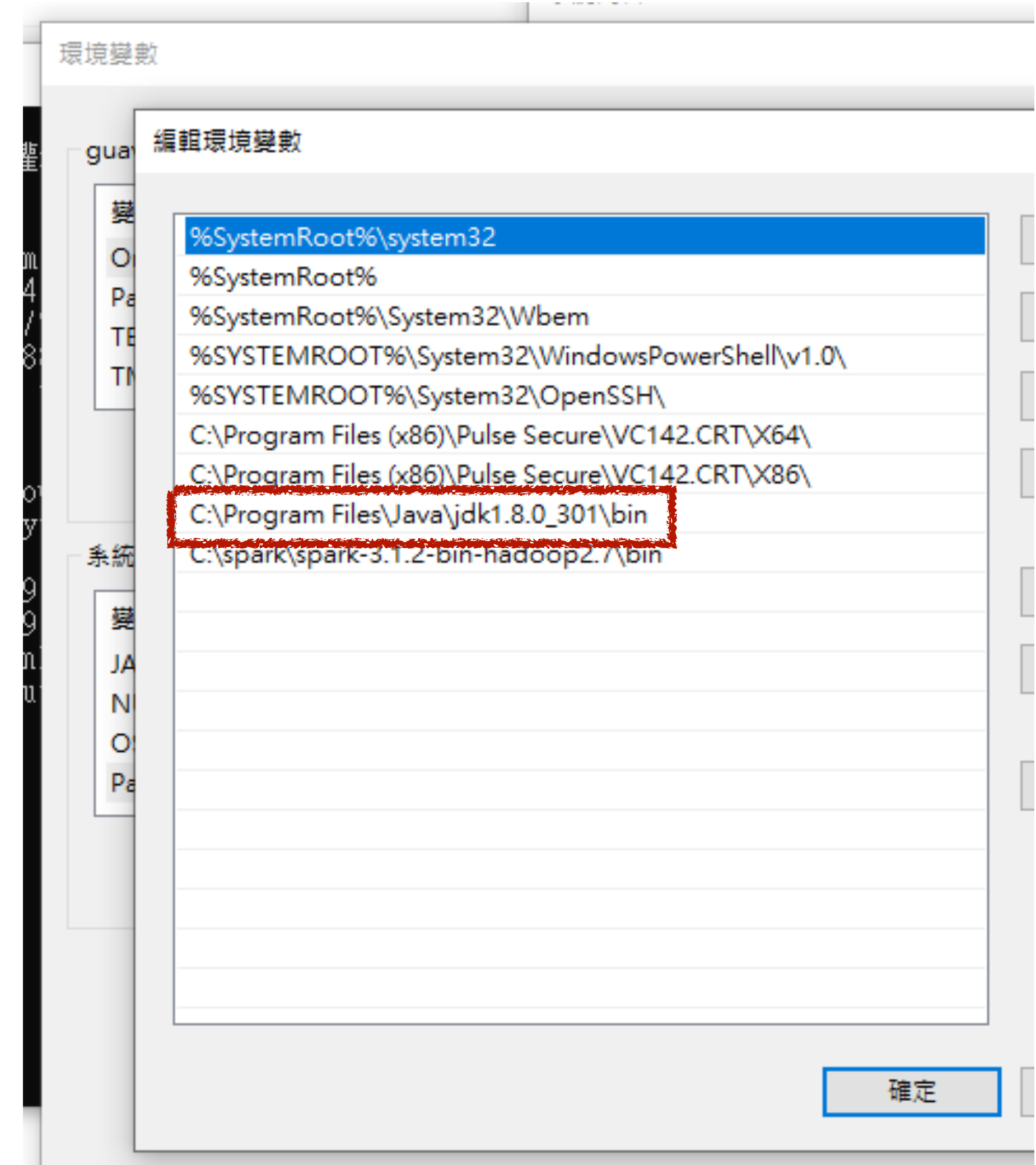- Fill in the **actual address** of the **bin** file in your jdk folder

# Step 3 - Spark Install and Setting
## Download Spark files

- Available download page: http://spark.apache.org/downloads.html

- Should select **Hadoop 2.7**

# Step 3 - Spark Install and Setting
## Setting environment variables

- At Environment Variables > System Variables

- Add **SPARK_HOME** and **HADOOP_HOME**

- Fill in the **actual address** of your spark folder
  (not the one in picture)
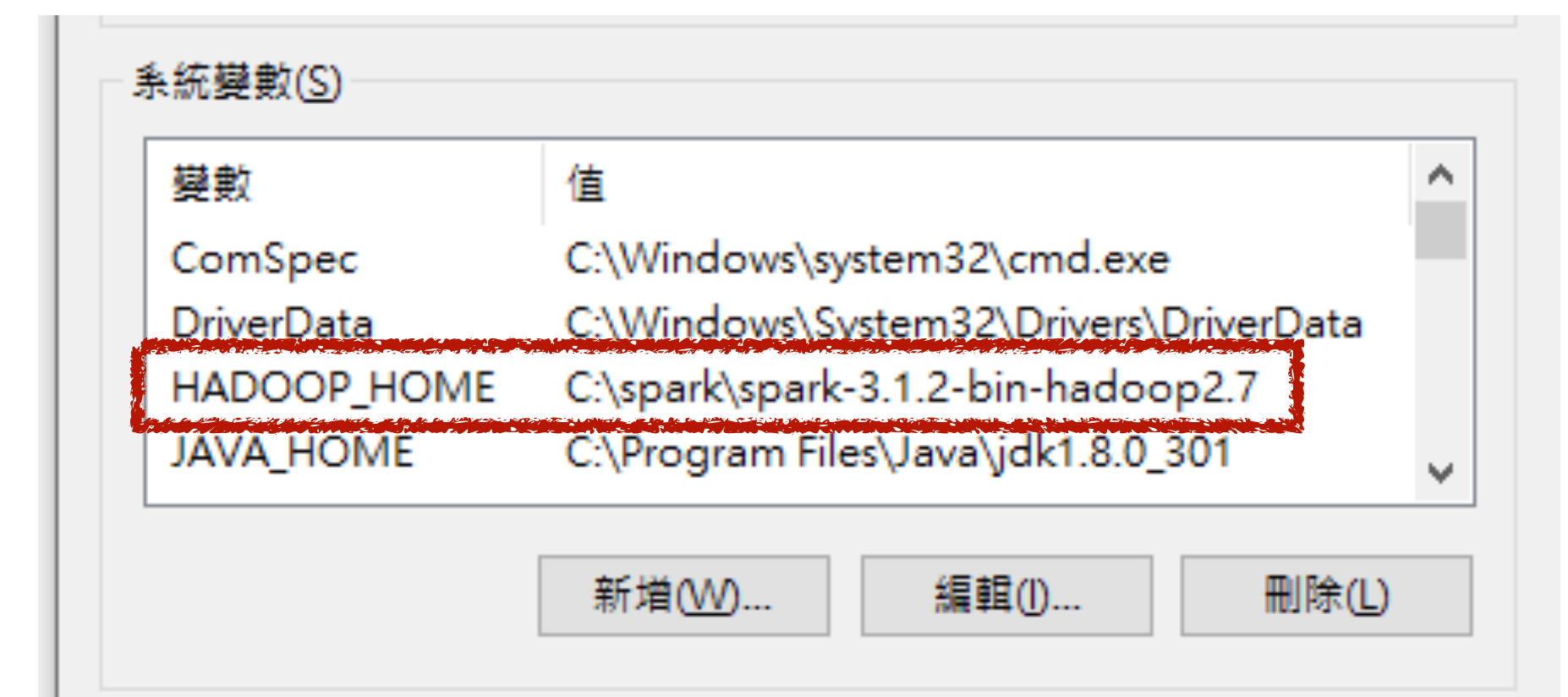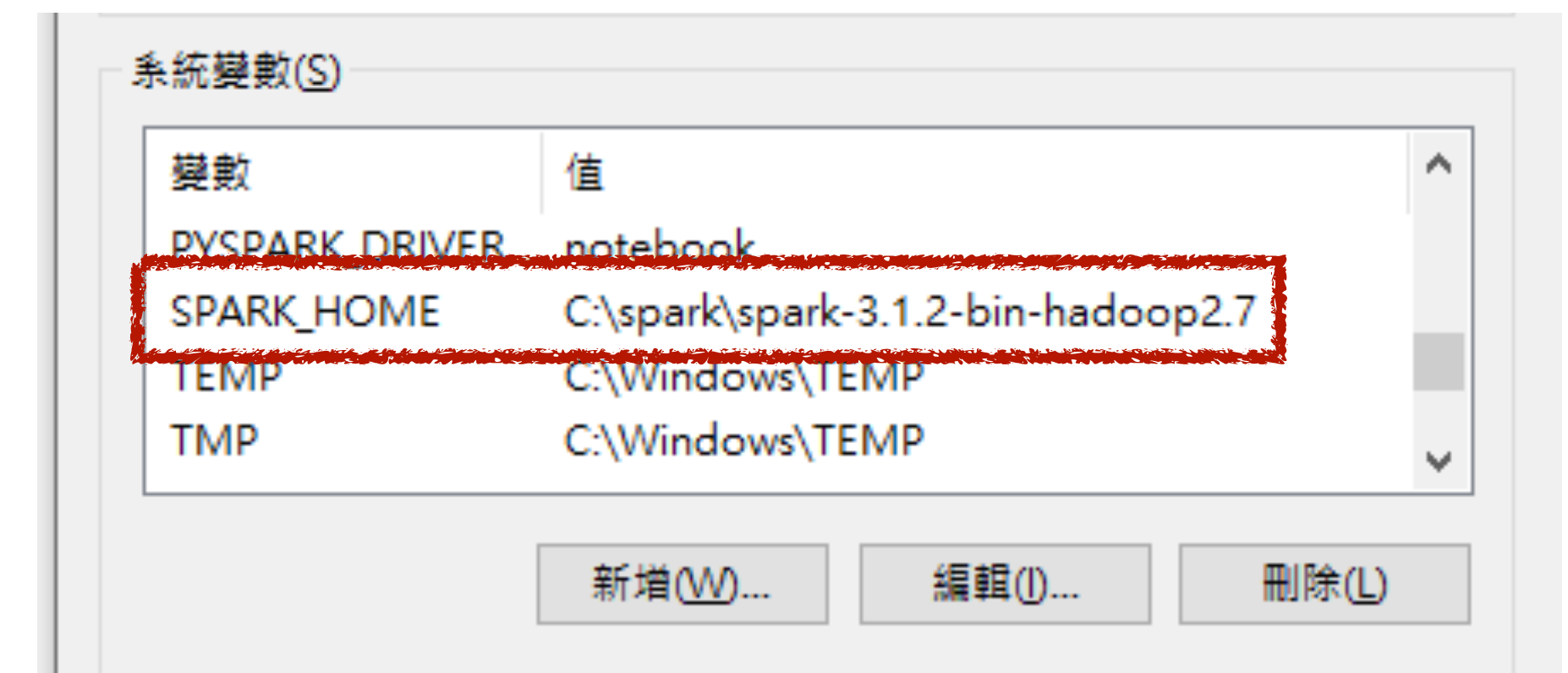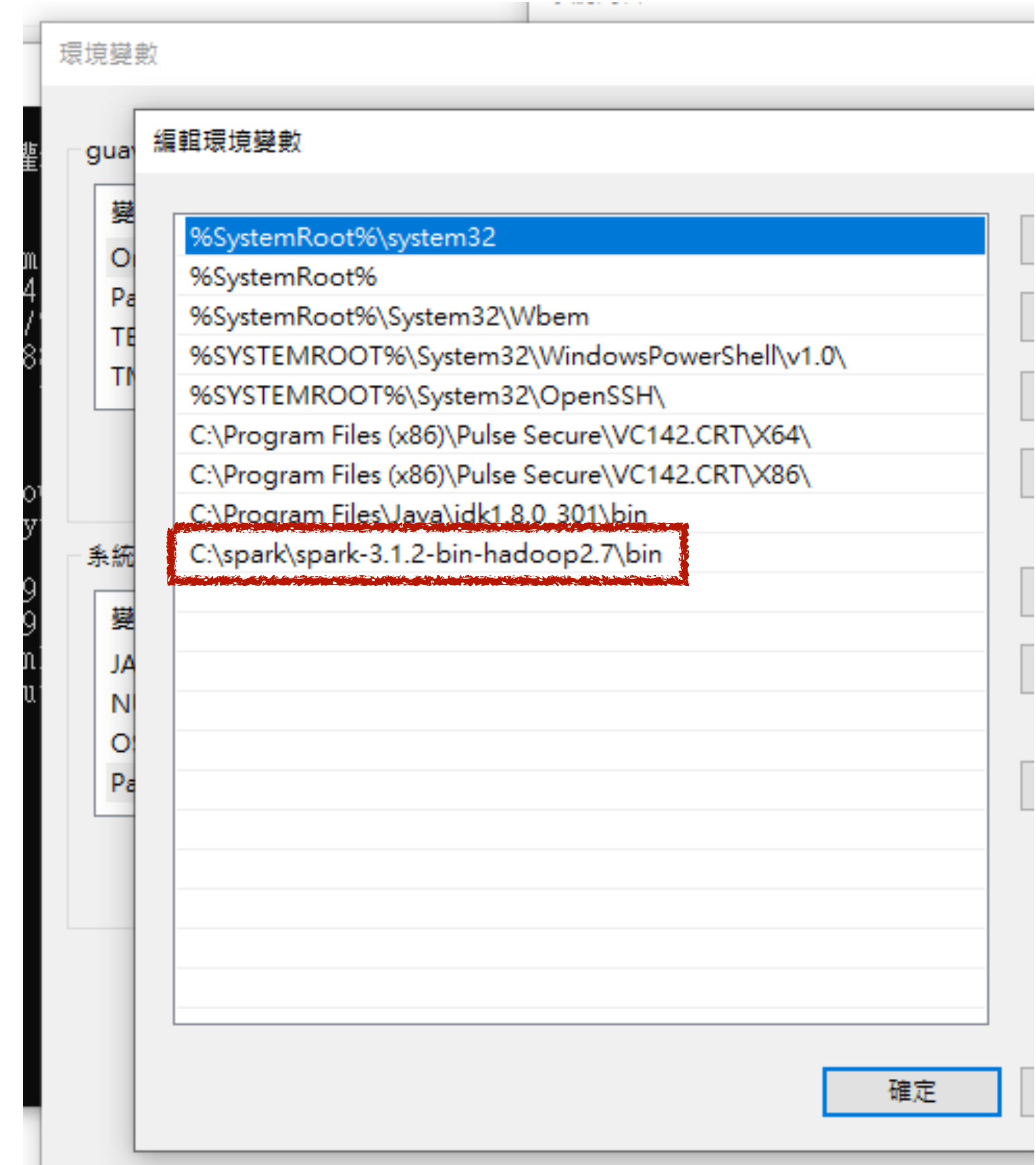
# Step 3 - Spark Install and Setting
## Setting environment variables

- At Environment Variables > System Variables

- Modify and add **PATH**

- Fill in the **actual address** of the **bin** file in your spark folder

# Step 4 - Hadoop winutils
## Download winutils.exe

- Available download page: https://github.com/steveloughran/winutils

- Download **winutils.exe** in hadoop.**2.7.1**

# Step 4 - Hadoop winutils
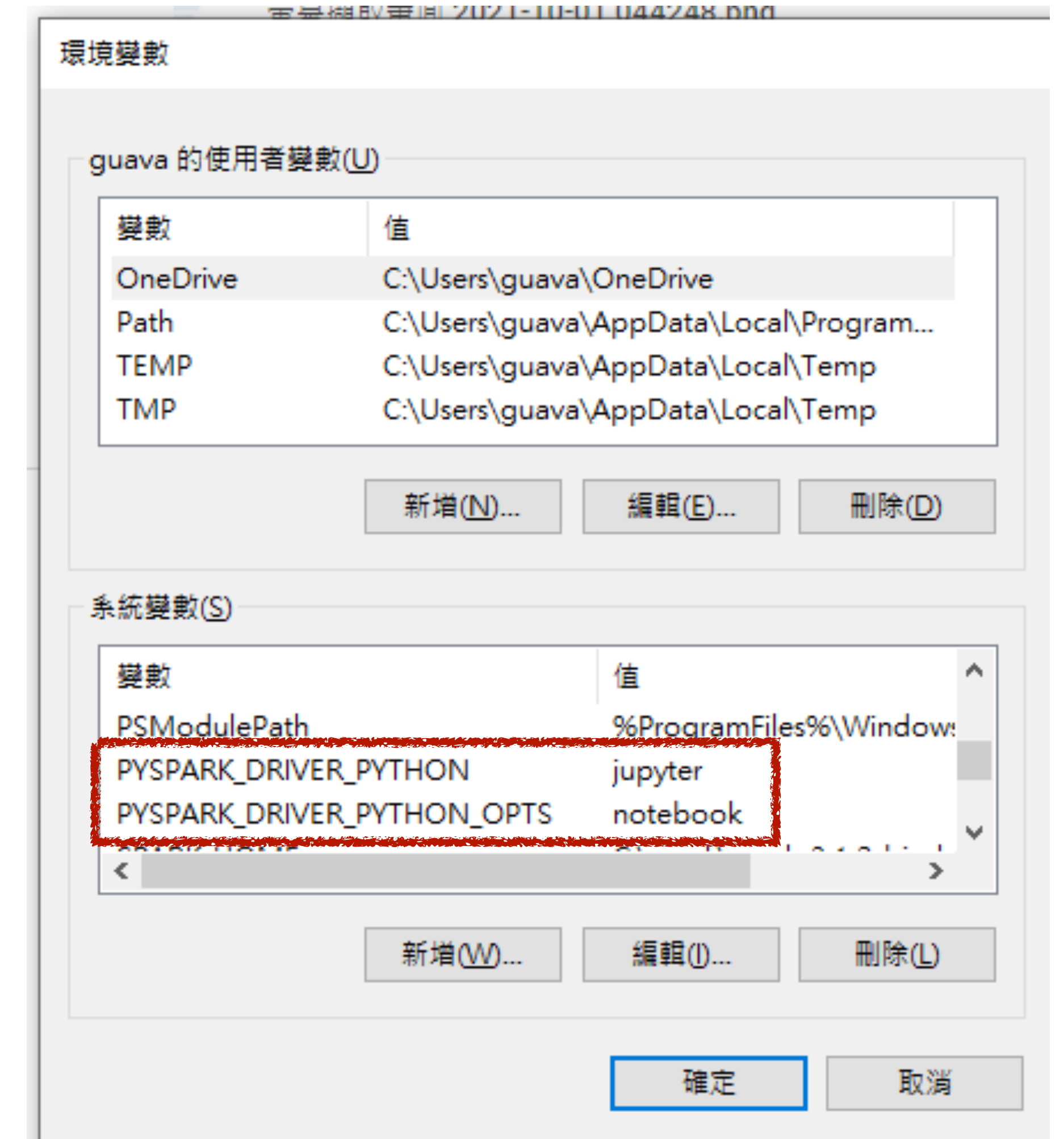## Set winutils.exe

– Put it into the bin folder in Spark

# Step 5 - jupyter notebook

## Install Jupyter and set environment variables

- Type  pip install jupyter  in cmd

- Add **PYSPARK_DRIVER_PYTHON**:
  **jupyter**

- Add **PYSPARK_DRIVER_PYTHON_OPTS**:
  **notebook**

# Step 6 - Testing

- Type `pyspark` in cmd

- The jupyter notebook page will open in browser

-

# Some do's and don'ts

- Avoid setting **Chinese words**, **space** or **special marks** in environment variables

- Make sure all addresses exist in your own device (Copy addresses on picture may not work)

- **Restart cmd** after updating environment variables