

LIVE:

[Multithreading] in Python for ML/AI  
& Multiprocessing]

Applied AI Course . com

## what we will cover

→ Multiprocessing & Threading  
for DataScience, ML &  
DL

→ code-walk-throughs

→ Introductory - Session

## what we will not cover

→ internals of OS [IPC, PM]

→ <https://gate.appliedcourse.com/course/5/operating-systems>

→ software engineering  
aspects

# Popular-Libraries-

## sklearn.linear\_model.LogisticRegression

```
class sklearn.linear_model.LogisticRegression(penalty='l2', *, dual=False, tol=0.0001, C=1.0, fit_intercept=True,
intercept_scaling=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto', verbose=0,
warm_start=False, n_jobs=None, l1_ratio=None)
```

[source]

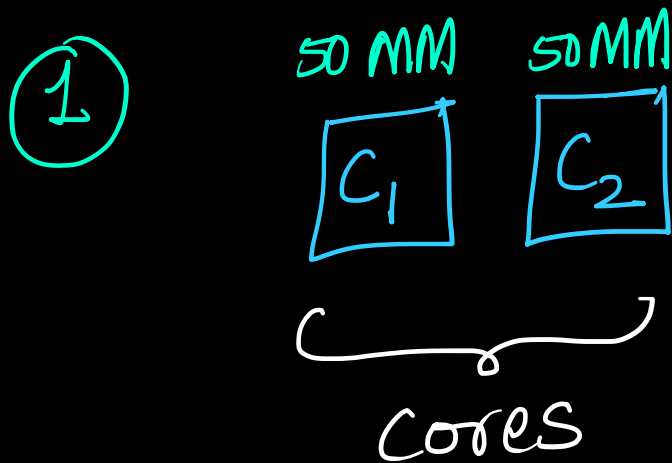
Logistic Regression (aka logit, MaxEnt) classifier.

→ joblib

Numpy & SciPy → BLAS

TensorFlow → multithreading, multi core, GPUs & distributed computing

Task: mean of 100 million numbers.

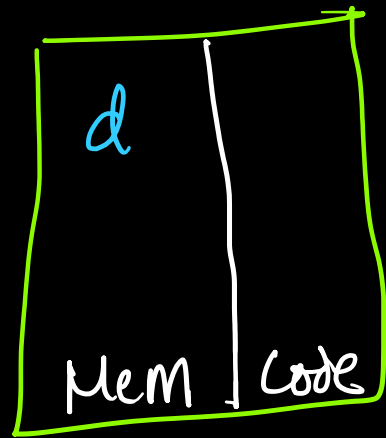


→  $m_1$  &  $m_2$

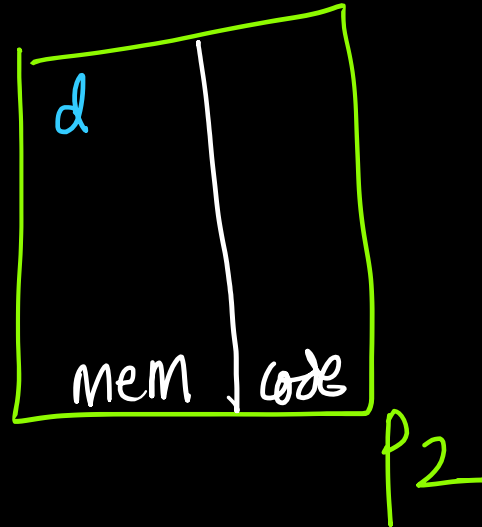
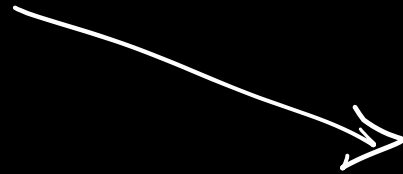
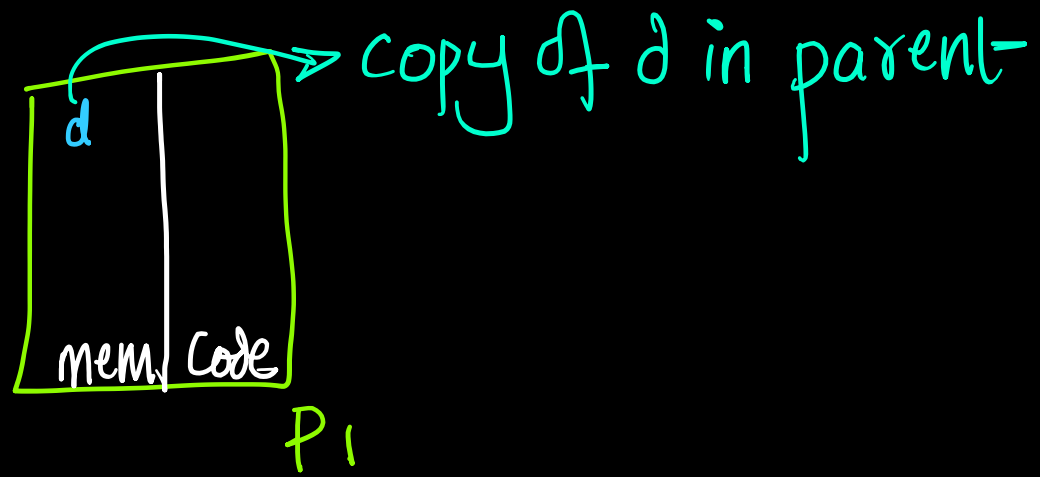
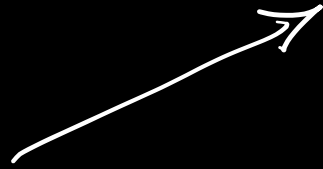
②

$$\frac{m_1 + m_2}{2} = \text{mean}$$

Multi-processing  
↓  
[code]



Parent

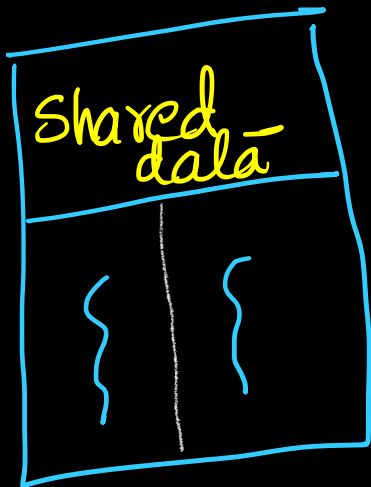


Tons of internal  
operating  
Systems

# Multithreading:

2 threads per core (Intel)

↑  
1 process per core @ a time



process

- lightweight
- faster context-switching

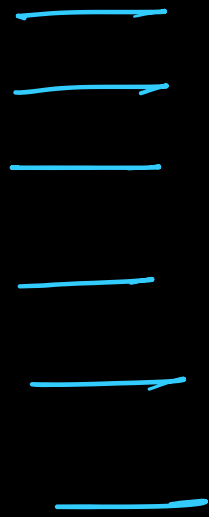
[code]

Combine multiprocessing & Threading



careful coding

Global Interpreter Lock → worst part of Python MT



Python  
Interpreter

$T_1$

$T_2$

python's memory-management is  
not thread-safe



Joblib

<https://joblib.readthedocs.io/en/latest/>

- Simple parallel-computing in Python
- Disk-caching of function outputs [code]
- Widely used : scikit-Learn

Joblib: Parallel [code]

→ only multi-processing

→ can also use multi-threading [GIT can slow down]

# Common Data science / ML / DL tasks for parallel processing

① Matrix & Vector products

② Data preprocessing

## ② Model-Training

- Logistic Regression

- Decision Trees

- Random Forests

- GBDT

### ③ Deep-Learning models

MLP

CNNs

Transformers

Parallelism for productionization [C/C++/Java]

Logistic Regression

GBDT / Random Forest

Deep Learning