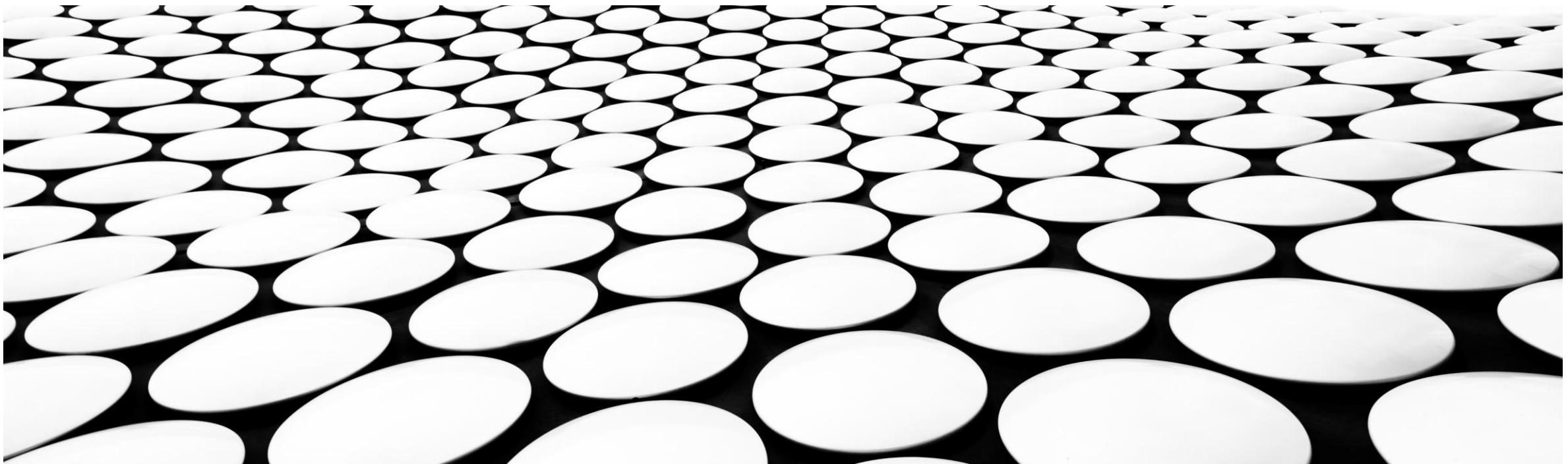

深度学习

邱怡轩



今天的主题

- 常用优化方法简介
- 深度学习计算环境



优化方法

传统模型

- 在传统统计模型中，经常使用牛顿法迭代
- 同时利用一阶导和二阶导的信息
- 如果参数数量是 p
- 一阶导是梯度 g ，大小为 $p \times 1$
- 二阶导是 Hessian 矩阵 H ，大小为 $p \times p$
- 牛顿法需计算 $H^{-1}g$ ，复杂度为 $O(p^3)$

一阶方法

- 对于深度学习模型，参数数量非常多
- 几乎只能依赖于梯度
- 通常称为一阶优化方法

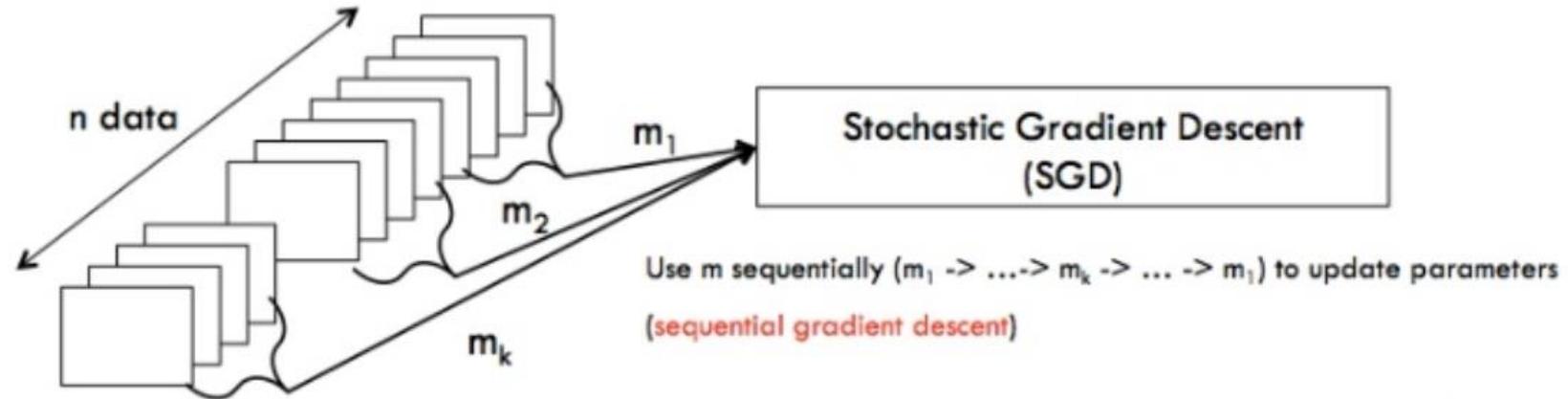
GD

- Gradient descent
- 梯度下降法，或最速下降法
- $\theta^{(k+1)} = \theta^{(k)} - \eta^{(k)} \cdot \nabla_{\theta} L(\theta^{(k)})$
- η 称为步长或学习率
- $\nabla_{\theta} L(\theta^{(k)})$ 是精确的梯度，由所有观测计算

SGD

- Stochastic gradient descent
 - 随机梯度下降
-
- $\theta^{(k+1)} = \theta^{(k)} - \eta^{(k)} \cdot \nabla_{\theta} l(\theta^{(k)})$
 - $\nabla_{\theta} l(\theta^{(k)})$ 是 $\nabla_{\theta} L(\theta^{(k)})$ 的一个无偏估计
 - 例如通过 mini-batch 计算而来

Mini-batch



思考题

- 数据共有 n 个观测，每个 mini-batch 大小为 m ，如有剩余则单独形成一个 mini-batch
- 如何简单计算 mini-batch 的数量？
- 用 Python 实现

SGD

- SGD 的重要意义在于，即使利用随机的、不精确的梯度
- 在一定条件下也可以保证收敛到导数为0的点
- $\sum_{k=1}^{\infty} \eta^{(k)} = \infty, \quad \sum_{k=1}^{\infty} (\eta^{(k)})^2 < \infty$
- 由统计学家 Herbert Robbins 和 Sutton Monro 在1951年提出
- A Stochastic Approximation Method

SGD

- $\eta^{(k)}$ 的选取
- 令 $S_n = \frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{n}$, $T_n = \frac{1}{1^2} + \frac{1}{2^2} + \dots + \frac{1}{n^2}$
- 则 $S_n - \log(n) \rightarrow 0.577216 \dots$, 故 $S_n \rightarrow \infty$
- $T_n \rightarrow \frac{\pi^2}{6} < \infty$

常用模板

```
1  obs_id = np.arange(n) # [0, 1, ..., n-1]
2  # Run the whole data set `nepoch` times
3  for i in range(nepoch):
4      # Shuffle observation IDs
5      np.random.shuffle(obs_id)
6
7      # Update on mini-batches
8  for j in range(0, n, batch_size):
9      # Create mini-batch
10     x_mini_batch = x[obs_id[j:(j + batch_size)]]
11     # Compute loss
12     loss = model(x_mini_batch)
13     # Compute gradient and update parameters
14     optimizer.zero_grad()
15     loss.backward()
16     optimizer.step()
```

改进 SGD

- SGD 虽然具有较好的理论性质
- 但在实际中会遇到各种挑战，如：
 - 确定合适的学习率 η ，过大导致优化不收敛，过小耗费大量迭代次数
 - 对每个参数使用了相同的学习率 η

改进SGD

改进 SGD

- SGD 虽然具有较好的理论性质
- 但在实际中会遇到各种挑战，如：
 - 确定合适的学习率 η ，过大导致优化不收敛，过小耗费大量迭代次数
 - 对每个参数使用了相同的学习率 η

Adagrad

- 核心思想是对每一个参数计算一个单独的 η
- 令 $g_i^{(k)}$ 为第 i 个参数在第 k 次迭代的导数
- SGD 即为 $\theta_i^{(k+1)} = \theta_i^{(k)} - \eta \cdot g_i^{(k)}$
- Adagrad 为 $\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\eta}{\sqrt{G_i^{(k)} + \varepsilon}} \cdot g_i^{(k)}$
- 其中 $G_i^{(k)} = (g_i^{(1)})^2 + \dots + (g_i^{(k)})^2$

Adagrad

- 优点：
 - 学习率自动衰减
 - 每个参数使用自适应的学习率
- 缺点：
 - 学习率衰减非常快，后期动力不足

RMSprop

- 对 Adagrad 加以改进
- 使用滑动平均计算 G_i
- $$G_i^{(k)} = \gamma G_i^{(k-1)} + (1 - \gamma) (g_i^{(k)})^2$$
- $$\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\eta}{\sqrt{G_i^{(k)} + \varepsilon}} \cdot g_i^{(k)}$$

Adadelta

- 与 RMSprop 独立地对 Adagrad 加以改进
- 加入分子的滑动平均，保持量纲一致
- $G_i^{(k)} = \gamma G_i^{(k-1)} + (1 - \gamma) (g_i^{(k)})^2$
- $\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\sqrt{D_i^{(k)} + \varepsilon} \cdot g_i^{(k)}}{\sqrt{G_i^{(k)} + \varepsilon}}$
- 分子为 $D_i^{(k)} = \gamma D_i^{(k-1)} + (1 - \gamma) (\Delta\theta_i^{(k)})^2$

Adam

- 对梯度也进行滑动平均
- $m_i^{(k)} = \beta_1 m_i^{(k-1)} + (1 - \beta_1) g_i^{(k)}$
- $v_i^{(k)} = \beta_2 v_i^{(k-1)} + (1 - \beta_2) (g_i^{(k)})^2$
- 修正偏差 $\hat{m}_i^{(k)} = m_i^{(k)} / (1 - \beta_1^k)$, $\hat{v}_i^{(k)} = v_i^{(k)} / (1 - \beta_2^k)$
- $\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\eta}{\sqrt{\hat{v}_i^{(k)}} + \epsilon} \cdot \hat{m}_i^{(k)}$

对比

- 理论性质的研究尚在进行
- 实际应用中往往比 SGD 更快收敛
- 前期用自适应的方法，后期用 SGD
- 见 [lec8-optimizer.ipynb](#)



Ali Rahimi 的演讲

<https://www.bilibili.com/video/BV1BW411Y78t>
13:27-16:00



计算环境

计算效率

- 神经网络的计算效率由众多不同的因素决定
- 算法收敛速度
- 软件实现
- 硬件设备
-

软件

- 主流的深度学习框架都对核心运算，如矩阵乘法、卷积等进行了高度优化

Tensorflow

- <https://www.tensorflow.org/>
- 免费、开源
- 主要由 Google 开发维护



PyTorch

- <https://pytorch.org/>
- 免费、开源
- 主要由 Facebook 开发维护



MXNet

- <https://mxnet.apache.org/>
- 免费、开源
- 由陈天奇、李沐等人发起
- 当前主要由 Apache 软件基金会和 Amazon 团队开发维护



MindSpore

- <https://www.mindspore.cn/>
- 免费、开源
- 由华为开发维护



对比

- 不同的软件框架之间并无绝对的优劣之分
- 语法通常非常相似
- 个人使用可依喜好选择
- 商业应用往往考虑兼容性和维护成本等因素

硬件

- 像绝大多数应用程序一样，神经网络模型可以运行在 CPU 上
- 随着深度学习的流行，更多专用设备如 GPU、TPU 等被用来加速计算

硬件

- GPU 的优势在于可以高度并行，特别适合神经网络的结构

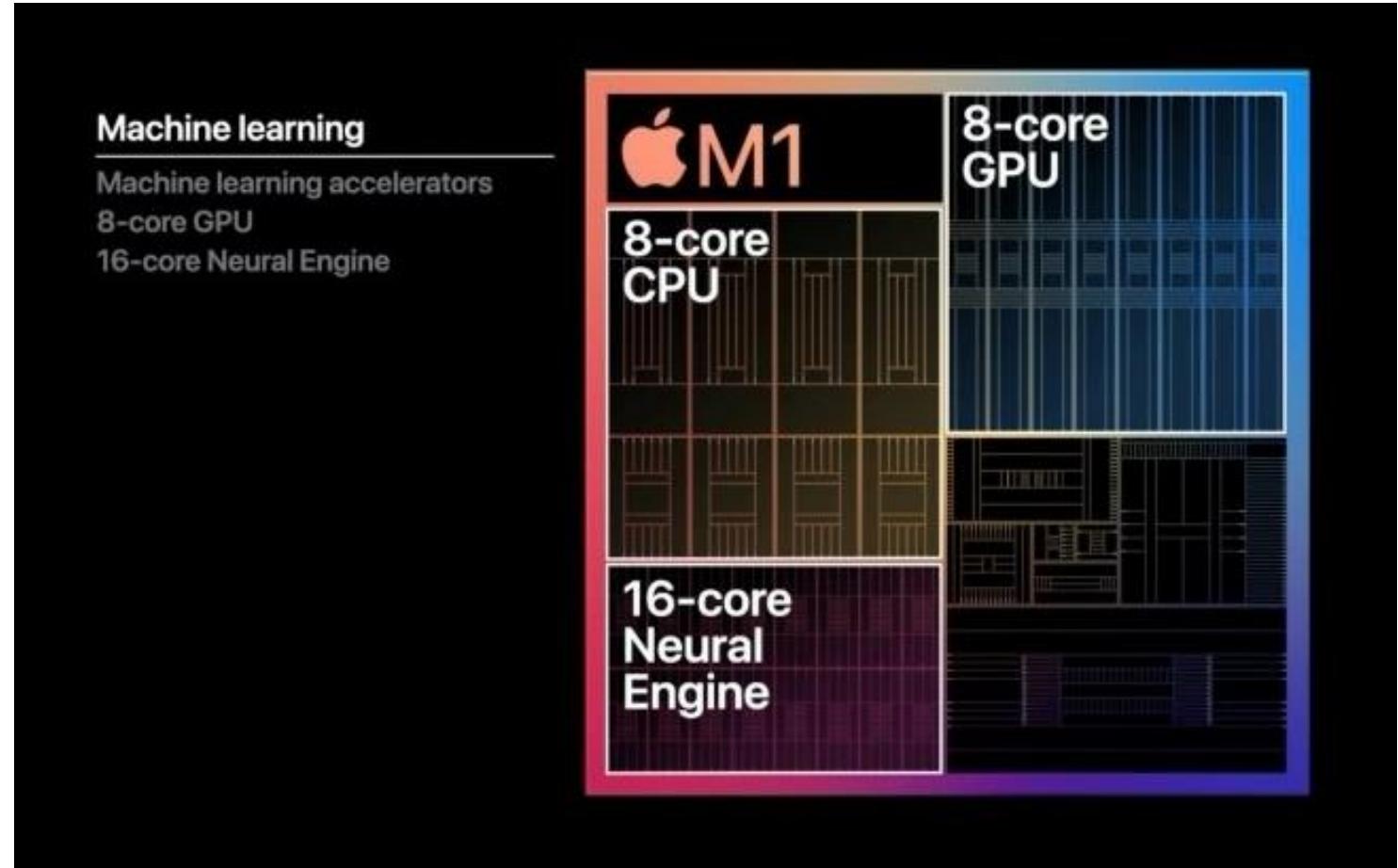
	GeForce RTX 5090	GeForce RTX 5080	GeForce RTX 5070 Ti
GPU Engine Specs:			
NVIDIA CUDA® Cores	21760	10752	8960
Shader Cores	Blackwell	Blackwell	Blackwell
Tensor Cores (AI)	5th Generation 3352 AI TOPS	5th Generation 1801 AI TOPS	5th Generation 1406 AI TOPS
Ray Tracing Cores	4th Generation 318 TFLOPS	4th Generation 171 TFLOPS	4th Generation 133 TFLOPS
Boost Clock (GHz)	2.41	2.62	2.45
Base Clock (GHz)	2.01	2.30	2.30

硬件

- 但并非 GPU 上运行的神经网络就一定比 CPU 上快
- 数据传输到 GPU 的计算核心需要时间
- GPU 单核的性能一般不如 CPU
- 数据量少、网络简单时并行效果不明显
- 对于较复杂的网络，GPU 的运算效率往往有很大的提升

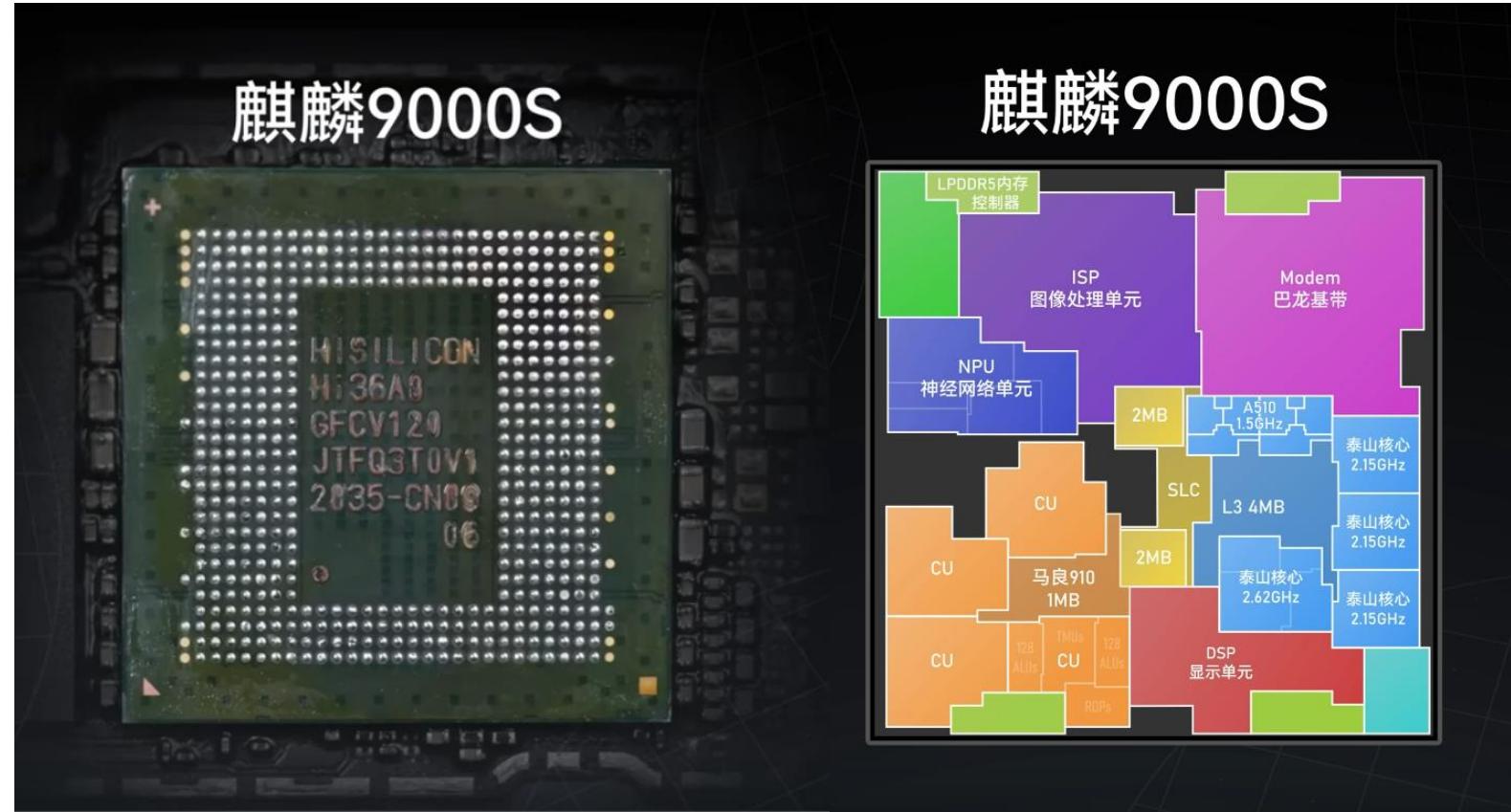
硬件

- 一些集成芯片上还带有专门的神经网络处理器



硬件

- 一些集成芯片上还带有专门的神经网络处理器



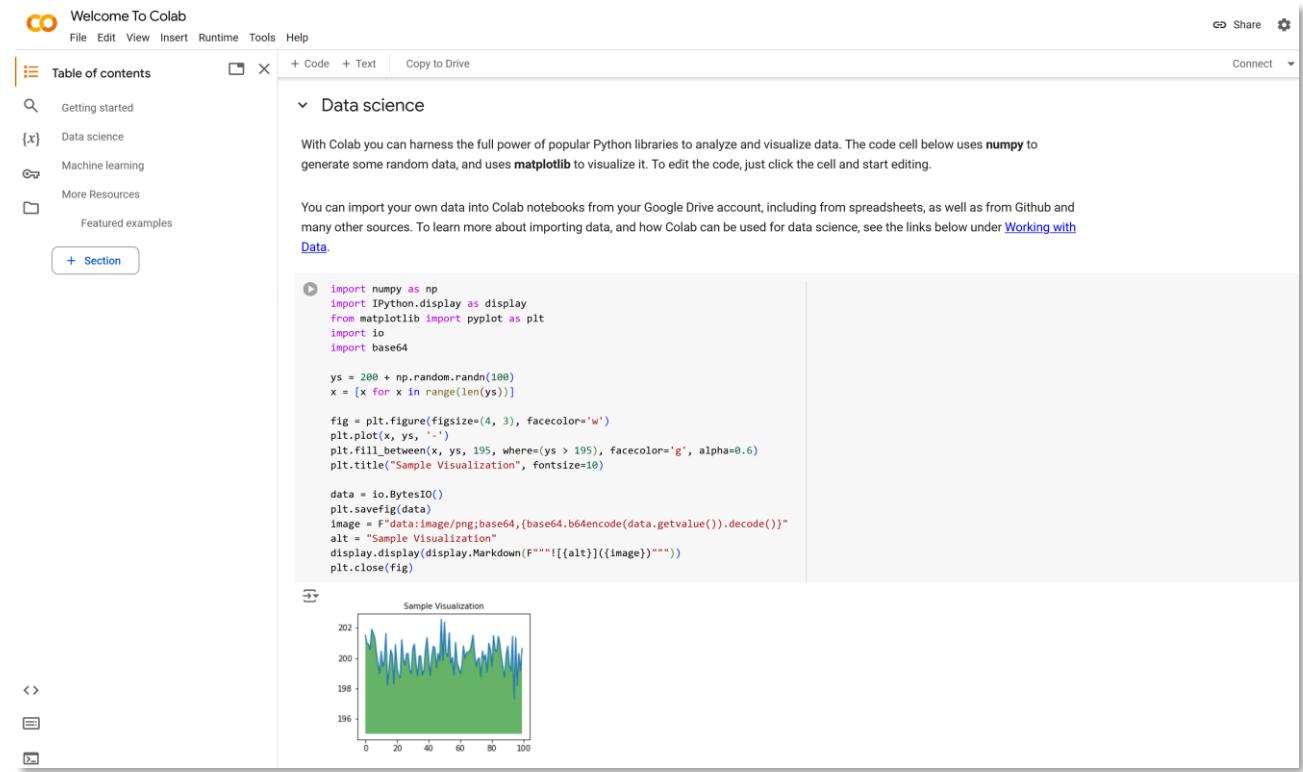
SUFE AI Lab

- <https://ai.sufe.edu.cn/>
- 上财校内计算资源
- 可使用 GPU

The image shows two screenshots of the SUFE AI Platform. The top screenshot is the homepage, featuring the university's logo and name, along with sections for Service Requests (including Home, Artificial Intelligence Case Center, Resource Application, CPU Service Application, GPU Service Application, My Resources, Help Documents, Plugin Download, and Insight), and a brief introduction to the AI platform. The bottom screenshot shows the 'Resource Application' section under 'Service Requests', specifically the 'CPU Service Application' tab, displaying various service options like Linux, Ubuntu, and Windows.

Google Colab

- <https://colab.research.google.com/>
- 免费计算资源（免费版有时间限制）
- 可使用 GPU



The screenshot shows the Google Colab interface. On the left, there's a sidebar with a 'Table of contents' section containing links to 'Getting started', 'Data science' (which is expanded), 'Machine learning', 'More Resources', and 'Featured examples'. Below this is a '+ Section' button. The main area has a header 'Welcome To Colab' with 'Share' and 'Connect' buttons. A code cell titled 'Data science' contains Python code for generating a plot. The code imports numpy, IPython.display, matplotlib.pyplot, io, and base64. It creates a random array 'ys', generates a plot with a green shaded area where ys > 195, and saves it as a base64 encoded image. The resulting plot is displayed below the code cell.

```
import numpy as np
import IPython.display as display
from matplotlib import pyplot as plt
import io
import base64

ys = 200 + np.random.randn(100)
x = [x for x in range(len(ys))]

fig = plt.figure(figsize=(4, 3), facecolor='w')
plt.plot(x, ys, '-')
plt.fill_between(x, ys, 195, where=(ys > 195), facecolor='g', alpha=0.6)
plt.title("Sample Visualization", fontsize=10)

data = io.BytesIO()
plt.savefig(data)
image = F"data:image/png;base64,{base64.b64encode(data.getvalue()).decode()}"
alt = "Sample Visualization"
display.display(display.Markdown(F""![[alt]]({{image}})"))

plt.close(fig)
```

Sample Visualization