# Lecture 09
# Push-button Deep Learning on the Cloud

H. Monajemi/DL. Donoho

Stats285, Stanford
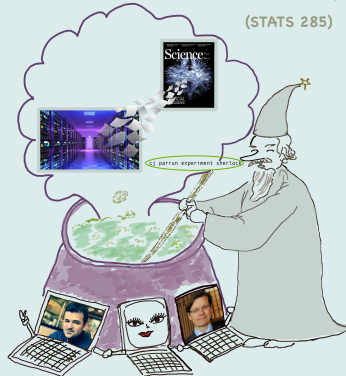
Nov/27/2017

**Stanford University**

# Stats 285 Fall 2017

## Disclaimer, 1: Google Images

*This document contains images obtained by routine Google Images searches. Some of these images may perhaps be copyright. They are included here for educational noncommercial purposes and are considered to be covered by the doctrine of* **Fair Use**. *In any event they are easily available from Google Images. It's not feasible to give full scholarly credit to the creators of these images. We hope they can be satisfied with the positive role they are playing in the educational process.*

**Stanford University**

# Outline

1. Computing Change

2. Deep Learning

3. Running Experiments on the Cloud

**Stanford University**

# *Research* computing demands are increasing
Shankar-Goodfellow Tweets

## *Research* computing demands are increasing
A personal story



Dear X,
I am trying to run GPU experiments on the cluster but these are all **pending** since all the CPUs are occupied. Can you please release 4 CPUs (for the 4 GPUs)? **Otherwise, we are stuck with no GPUs.**

**Stanford University**

## *Research* computing demands are increasing
A personal story

**X**

I changed the partition for the pending jobs. **Let's submit your jobs tonight** so that I can bring back those jobs tomorrow early morning.



But I need to **experiment and debug my code** ... how can I submit "my jobs tonight?"

**Stanford University**

## Stats285 theme: Explosion of Computing Resources

Cloud Paradigm:

- Billions of smart devices each drive queries to cloud servers
- Millions of business relying on cloud for all needs

Symbiosis of cloud and economy is *lasting* and *disruptive*.

Cloud provides *any user* **same-day** delivery:

- Tens to hundreds of thousands of hours of CPU
- Pennies per CPU hour
- $\approx 50$ cents per GPU hour

Any user can consume *1 Million CPU hours* over a few days for a few \$10K's.

Stanford University

# Cloud providers offer many services

| Google Cloud Platform | Amazon Web Services[7] | Microsoft Azure[8] |
|---|---|---|
| Google Compute Engine | Amazon EC2 | Azure Virtual Machines |
| Google App Engine | AWS Elastic Beanstalk | Azure Cloud Services |
| Google Container Engine | Amazon EC2 Container Service | Azure Container Service |
| Google Cloud Bigtable | Amazon DynamoDB | Azure Cosmos DB |
| Google BigQuery | Amazon Redshift | Microsoft Azure SQL Database |
| Google Cloud Functions | Amazon Lambda | Azure Functions |
| Google Cloud Datastore | Amazon DynamoDB | Cosmos DB |
| Google Storage | Amazon S3 | Azure Blob Storage |

source: Wikipedia

Stanford University

# Cloud can accommodate research computing needs
## Papyan's case



Why frustrated? create your own GPU cluster on the cloud

45 min later...



Hi X,
**I created a cluster following homework2 in stats285**. So I have computational resources now.

**Stanford University**

# Cloud can accommodate research computing needs
Shankar's case

## Computing Change is real!

- Computational Science/Computation-enabled discovery is becoming mainstream!
- PhD Students are expected to conduct **1 million CPU-hour** of computation
- Personal Laptops $\rightarrow$ Shared Clusters $\rightarrow$ Personal Clusters

**Stanford University**

## Adapting to Computing Change

- Just as Climate Change demands adaptation,
- Computing Change demands adaptation:
    - **Pose** bold research **hypotheses** to settle computationally
    - **Design massive computing experiments**
    - **Adopt** painless computing **frameworks**
    - **Raise money** to pay for cloud-based computing
    - *Push a button*

**Stanford University**

How to get rid of the computing interface pain?

# We need to rethink the way we do computational experiments

**Stanford University**

## What does an experiment involve?

In our telling, a computational experiment involves:

1. **Precise Specification** (define metric and parameters)
2. **Execution and management** of all the jobs
3. **Harvesting** of all the data generated by all the jobs
4. **Analysis** of the data
5. **Reporting** of results.

Today we add to this list: **Building** of compute clusters

The painless computing paradigm should seamlessly integrate and automate all these tasks

**Stanford University**

# Many open-source frameworks offer automation
painless (push-button) massive computing

- Building Cloud Clusters
    - **ElastiCluster** (Riccardo Murri)
    - StarCluster (MIT)
- Experiment Management Systems (Laptop $\rightarrow$ Cluster)
    - **CJ** (Yours Truly)
    - CodaLab (Percy Liang)
    - PyWren(serverless) (Eric Jonas)
- Machine Learning and Statistics
    - **PyTorch**, Tensorflow, CNTK, Theano, Keras
    - Spark, Dask

**Stanford University**

# Today's focus
Push-button Massive Computational Experiments

# Outline

**Stanford University**

## Classical Statistics: Linear relationship

Given $n$ realizations $(x_i, y_i)$,

$$y_i = \beta^T x_i + \epsilon_i$$

$$x_i, \beta \in \mathbb{R}^p, \quad y_i \in \mathbb{R}, \quad i = 1, \ldots, n$$

**Stanford University**

## Non-linear relationship

Given $n$ realizations $(x_i, y_i)$,

$$y_i = \Phi(x_i; \Omega) + \epsilon_i$$

- How to choose $\Phi$?
- An example of a classical approach:
  - $\Phi(x; c) = \sum_\ell c_\ell K(x_\ell, x)$ with known *kernel* $K$
  - Find $c_\ell$ such that data is reproduced by the model

**Stanford University**

## Neural Nets Approach

$$y_i = \Phi(x_i; \Omega) + \epsilon_i$$

$$\Phi(x; W, \beta) = \beta^T \sigma(W^T x)$$

- $W \in \mathbb{R}^{p \times d_1}$, $\beta \in \mathbb{R}^{d_1}$
- $\sigma$: A nonlinearity typically
    - *non-negative soft-thresholding (ReLU)*
    - *logistic function (sigmoid)*
- Find $(W, \beta)$ such that data is reproduced (a.k.a **Training**).

**Stanford University**

## Visual illustration

$$i\text{-th row: } \sigma(W^T x)_i \equiv h_i$$



- Perceptron, the basic block (Rosenblatt, 1957)

## Visual illustration

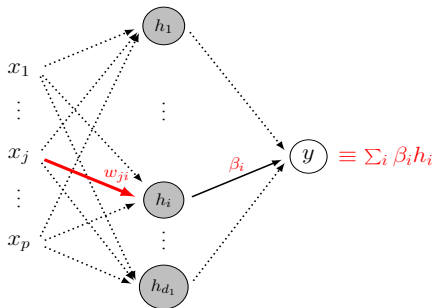$$h_i = \sigma(W^T x)_i, \quad i = 1, \ldots, d_1$$



- Single-Layer Perceptron

## Visual illustration

$$\Phi(x; W, \beta) = \beta^T \sigma(W^T x)$$



- Univariate output

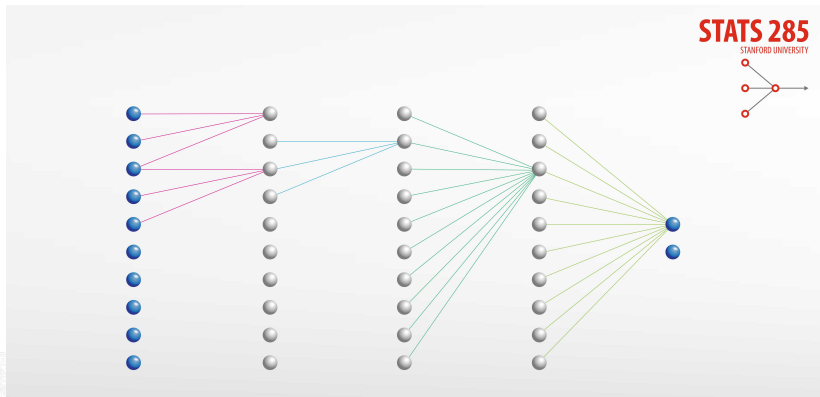## Deep Neural Nets

$$y_i = \Phi(x_i; \Omega) + \epsilon_i$$

$$\Phi(x; W_1, \ldots, W_L) = W_L^T \sigma\Big(W_{L-1}^T \ldots \sigma(W_1^T x)\Big)$$

- $W_\ell \in \mathbb{R}^{d_{\ell-1} \times d_\ell}, \quad d_0 = p$
- **Training**: find $(W_1 \ldots, W_L)$.
- highly over-parametrized (# unknowns $\gg n$ )
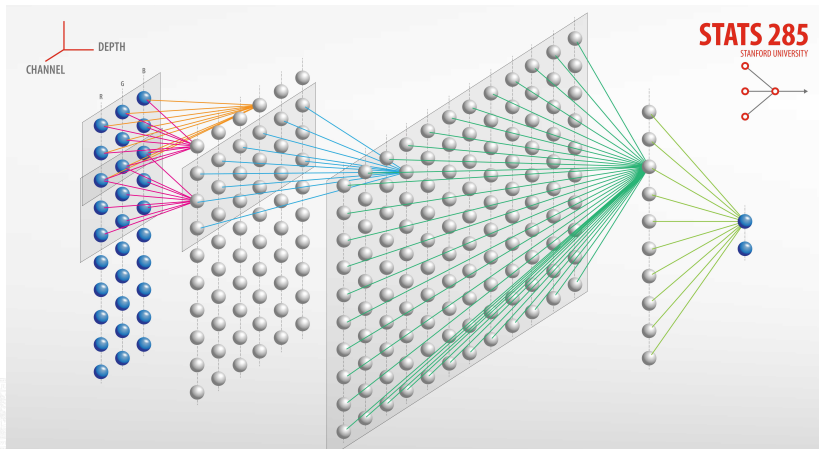
**Stanford University**

# Visual illustration, Deep Nets - 2D

- Locality (convolution)
- Weight Sharing

# Visual illustration, Deep Nets - 3D

- Locality (convolution)
- Weight Sharing

# Back-propagation – derivation
derivation from LeCun et al. 1988

Given $n$ training examples $(x_i, y_i) \equiv$ (input,target) and $L$ layers

- Constrained optimization

$$\min_{W,x} \qquad \sum_{i=1}^{n} \|h_i(L) - y_i\|_2$$

$$\text{subject to} \quad h_i(\ell) = \sigma_\ell \Big[ W_\ell h_i\left(\ell - 1\right) \Big],$$

$$i = 1, \ldots, n, \quad \ell = 1, \ldots, L, \; h_i(0) = x_i$$

- Lagrangian formulation (Unconstrained)

$$\min_{W,x,B} \mathcal{L}(W, x, B)$$

$$\mathcal{L}(W, x, B) = \sum_{i=1}^{n} \left\{ \|h_i(L) - y_i\|_2^2 + \right.$$

$$\left. \sum_{\ell=1}^{L} B_i(\ell)^T \left( h_i(\ell) - \sigma_\ell \Big[ W_\ell h_i\left(\ell - 1\right) \Big] \right) \right\}$$

## Back-propagation – derivation

### Forward pass, $\frac{\partial \mathcal{L}}{\partial B}$

$$h_i(\ell) = \sigma_\ell \Big[ \underbrace{W_\ell h_i(\ell-1)}_{A_i(\ell)} \Big] \quad \ell = 1, \ldots, L, \quad i = 1, \ldots, n$$

### Backward (adjoint) pass, $\frac{\partial \mathcal{L}}{\partial h}$, $z_\ell = [\nabla \sigma_\ell] B(\ell)$

$$z(L) = 2\nabla \sigma_L \Big[ A_i(L) \Big] (y_i - h_i(L))$$

$$z_i(\ell) = \nabla \sigma_\ell \Big[ A_i(\ell) \Big] W_{\ell+1}^T z_i(\ell+1) \quad \ell = 0, \ldots, L-1$$

### Parameter update, $W \leftarrow W + \lambda \frac{\partial \mathcal{L}}{\partial W}$

$$W_\ell \leftarrow W_\ell + \lambda \sum_{i=1}^n z_i(\ell) h_i^T(\ell-1)$$

## Back-prop in PyTorch

- Forward pass

  ```
  output    = model(input)
  loss      = loss_fn(outputs, target)
  ```

- Parameter Update

  ```
  loss.backward()    # computes
  optimizer.step()   # adds to
  ```
  $\sum_i z_i h_i^T$

  $W^k$

**Stanford University**

# Deep Learning Experiments

**1** ETL

```
trainloader, testloader = dl.torch.data.etl(dataset=data_set, ...)
```

**2** Define model (network)

```
model = dl.torch.models.mini_alexnet()
```

**3** Define loss and optimizer

```
loss_fn     = torch.nn.CrossEntropyLoss()
optimizer   = torch.optim.SGD(model.parameters(), ...)
```

**4** Train

```
dl.torch.train(model, loss_fn, optimizer, trainloader, testloader, ...)
```

**Stanford University**

# Understanding deep learning requires rethinking generalization (Zhang et al.)

# Example of Replication

90 epochs in less than 5 min on Google Compute Engine

# Outline

**Stanford University**

# Stats285 on the cloud!

# Push-button Massive Computational Experiments

## *Literally* Push-button!

Action items for conducting experiments:

1) **elasticluster** start **gce**

2) **cj** parrun train.py **gce** -alloc "--gres=gpu:1"

**Stanford University**

# Build your cluster
## `elasticluster start gce`

- Setup Google Cloud Billing (*Yes, they will charge you!*)
  1. `gmail username`
  2. `project_id`
  3. `client_id`
  4. `client_secret`
- Install Docker on your machine
  - removes the pain and ir-reproducibility of software install
- Create your cluster using dockerized ElastiCluster

```
docker image pull stats285/elasticluster-gpu

docker run -it stats285/elasticluster-gpu

elasticluster start gce
```

Stanford University

# Fire and forget!
## cj parrun train.py gce

- add your cluster info to `CJ_install/ssh_config`

```
[gce]
Host        35.199.171.137
User        hatefmonajemi
Bqs         SLURM
Repo        /home/hatefmonajemi/CJRepo_Remote
Python      python3.4
Pythonlib   pytorch:cuda80:-c soumith
[gce]
```

- Fire up your jobs and track them using **ClusterJob**

```
cj parrun train.py gce -alloc "--gres=gpu:1"
```

**Stanford University**

# Your assignment is online

## Massive Computational Experiments, Painlessly (STATS 285)

Stanford University, Fall 2017

### Assignment 02

In this assignment, we will conduct a collaborative project testing certain theoretical hypotheses in Deep Learning. In particular, each of you will build **your own personal SLURM cluster** on Google Compute Engine (GCE) using elasticluster and then run massive computational experiments using clusterjob. We then collect and analyse all the results you will generate and document our observations. Please follow the following step to setup your cluster and run experiments. This documents only contains the detail of setting up your cluster and testing that it works properly with
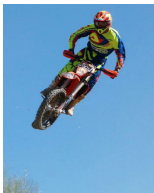
Stanford University

## Summary

- Small computations $\rightarrow$ Massive computations
- Interactive model (copy-paste) $\rightarrow$ experiment management systems (EMS)
- Personal Laptops $\rightarrow$ Shared Clusters $\rightarrow$ Personal Clusters
- Expansion by factors of 1000's in immediate computing capacity
- Rise of frameworks takes away pain of massive computing
- Deep Learning is now a technology (Thanks to frameworks)
- Lots of opportunities for semi-empirical/semi-theoretical Deep Learning studies

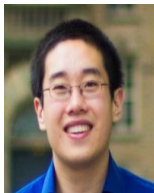**Stanford University**

## acknowledgements

- People



J. Lozano     X.Y. Han     R. Murri     V. Papyan

- Google Cloud Platform Education Grants
- ElastiCluster Team

**Stanford University**