

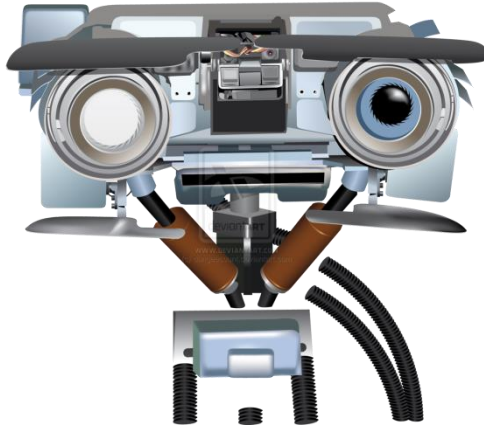
Announcements

Reminder: Class challenge out! Ends December 10th

Assignment Project Exam Help

- Lab this week: final review, class survey
<https://powcoder.com>

Add WeChat powcoder



Assignment Project Exam Help

On-device Machine Learning

<https://powcoder.com>

Add WeChat powcoder

CS542 Fall 2020

Creating an AI Model

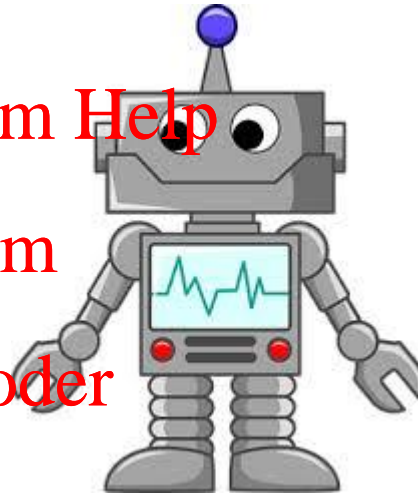


model development

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



model deployment

What if the model you trained is too big/takes too long to run?

Reminder: memory representation

- Int – 4 bytes

- Float – 4 bytes

Assignment Project Exam Help

<https://powcoder.com>

- Double – 8 bytes

Add WeChat powcoder

- Binary variable – 1 byte

Aha! We can reduce memory by $1/4^{\text{th}}$ if we can binarize the weights/features!

Binary features enable us to use hamming distance for fast similarity computation!

Why should I care?

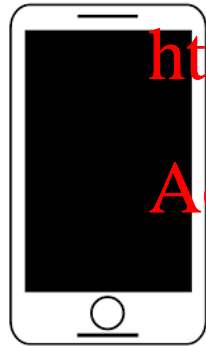
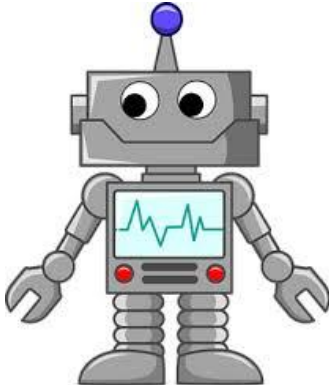
resource scarce systems

efficient features for search

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Goals

- Efficiency
 - Inference speed
 - Memory reduction (either RAM or disk space)
 - Limit Performance loss
 - ~~Training time?~~
- Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

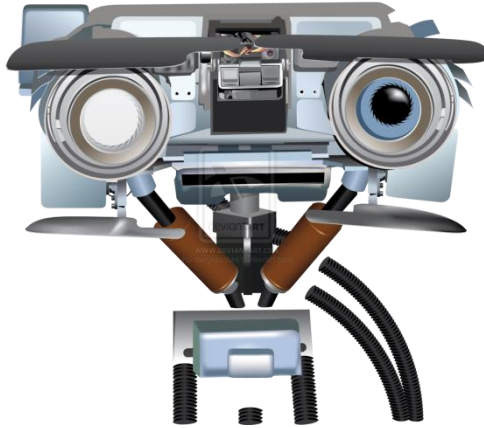
Today

- Network/Feature Quantization
- Parameter Pruning
- Knowledge Distillation

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Assignment Project Exam Help

<https://powcoder.com>

Feature Quantization

Add WeChat powcoder

CS542 Fall 2020

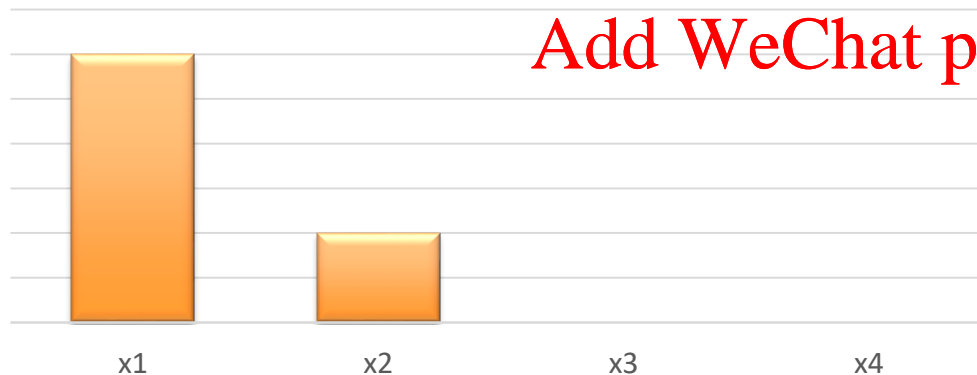
Simple feature quantization method - threshold

- Find the mean μ of each feature x_i in the training set, then binarize using $\mu > x_i$

Assignment Project Exam Help

<https://powcoder.com>

variance per feature

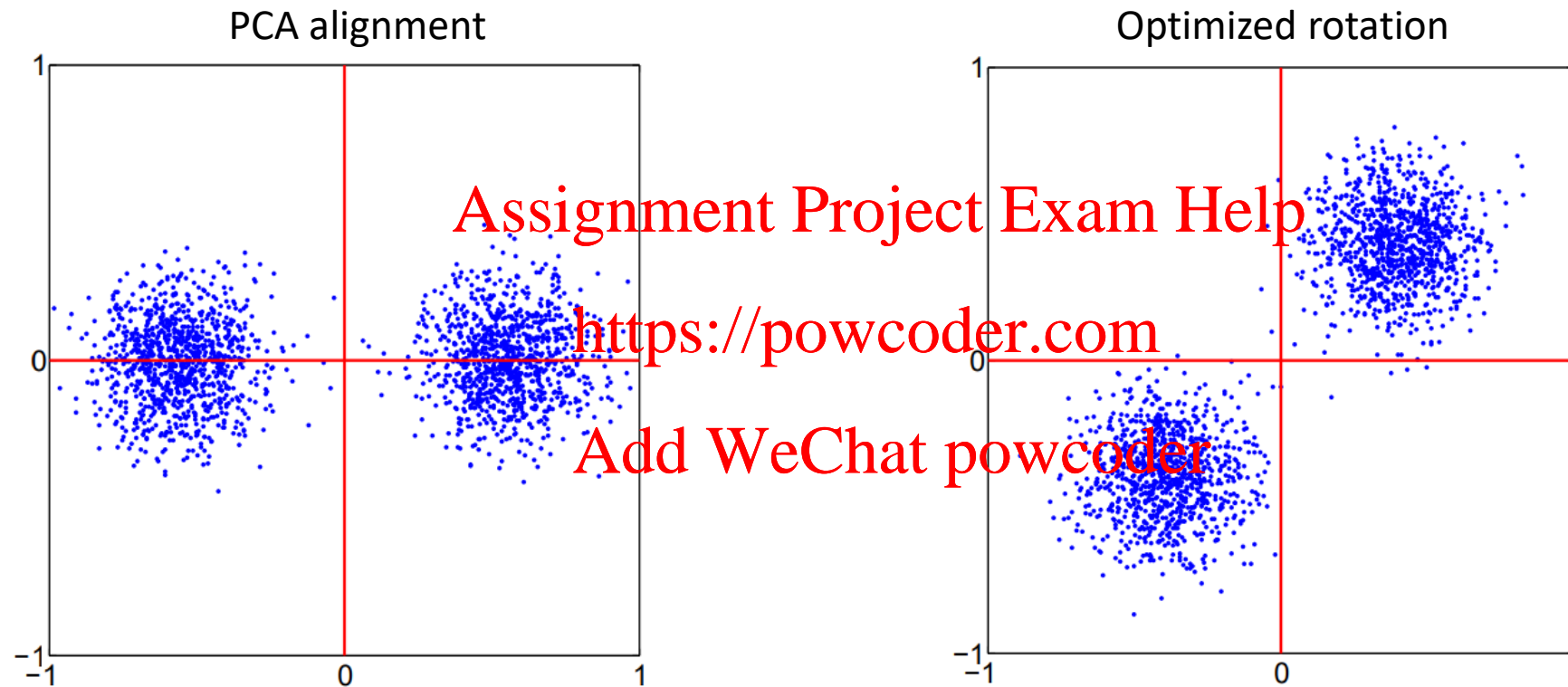


Add WeChat powcoder

Some features have more information than others

Also assumes the initial feature representation size is small enough

Iterative quantization



Paper reference: [Iterative Quantization: A Procrustean Approach to Learning Binary Codes](#). Y. Gong and S. Lazebnik, CVPR 2011.

Quantization optimization

For orthogonal rotation matrix R and projected features V , minimize:

Assignment Project Exam Help

$$Q(B, R) = \|B - VR\|_F^2$$

<https://powcoder.com>

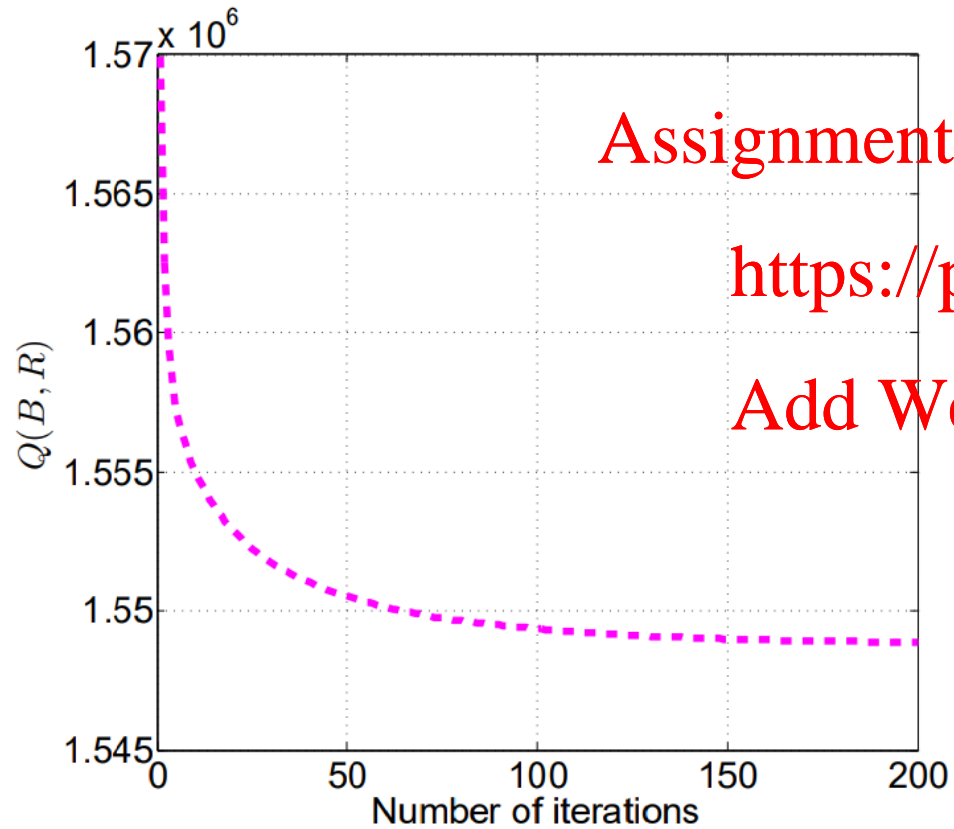
where B is a binary code matrix. Add WeChat powcoder

But it is dependent on two variables!

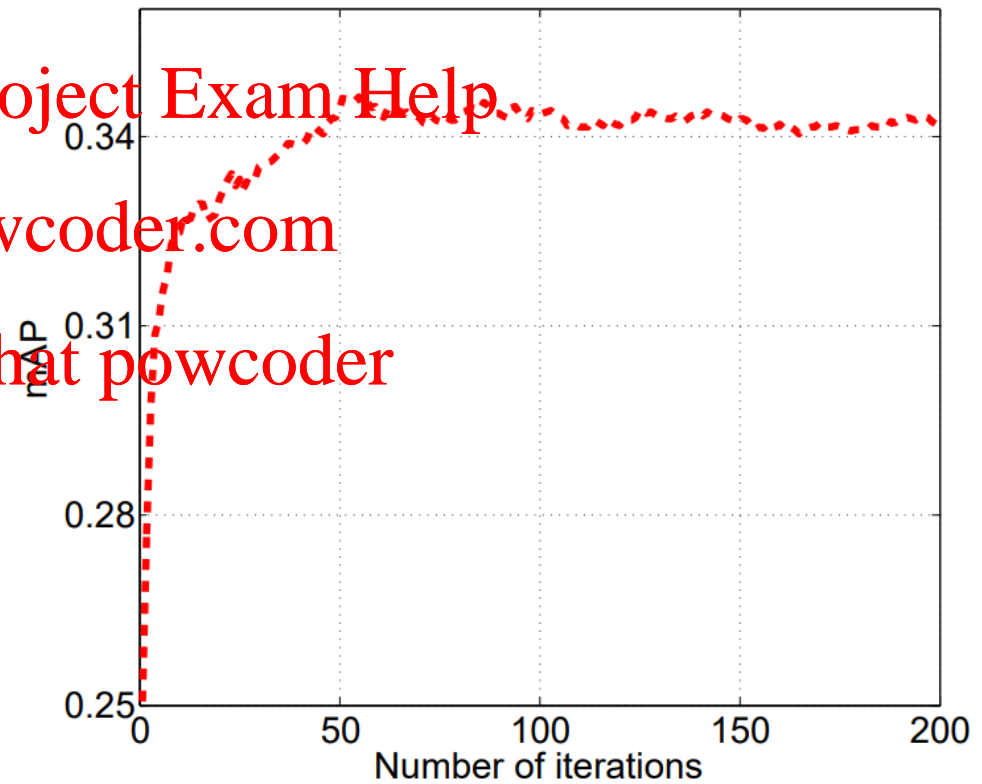
Fix one, update the other

Performance over time

Cost function value



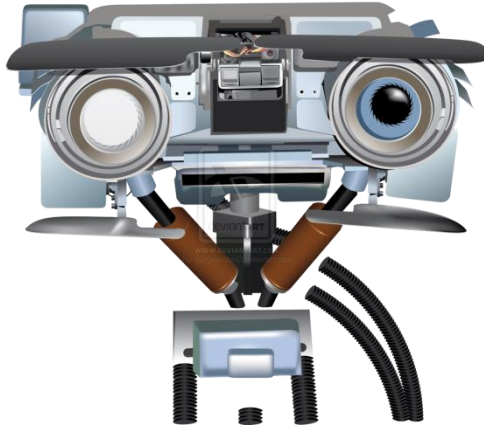
Euclidean neighbor precision



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Assignment Project Exam Help

<https://powcoder.com>

Network Quantization

Add WeChat powcoder

CS542 Fall 2020

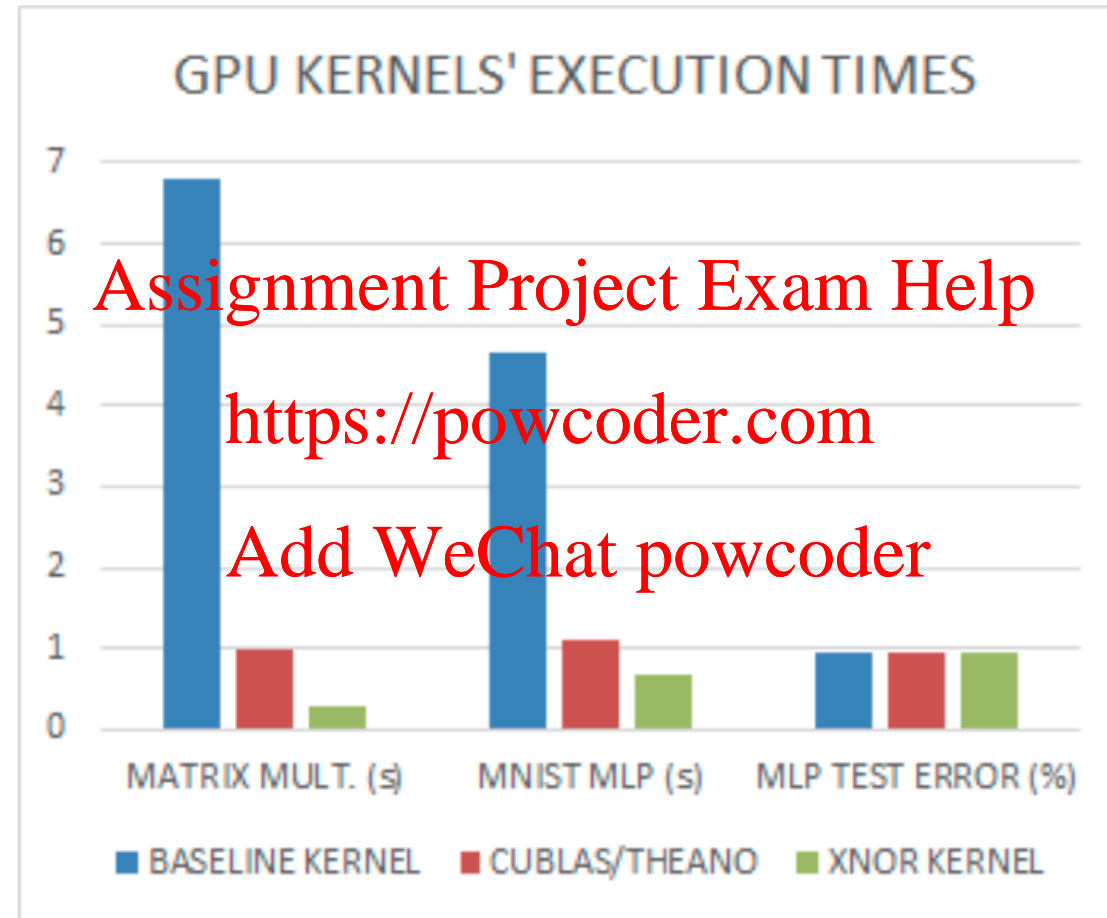
Convolutions in binary networks

Assignment Project Exam Help

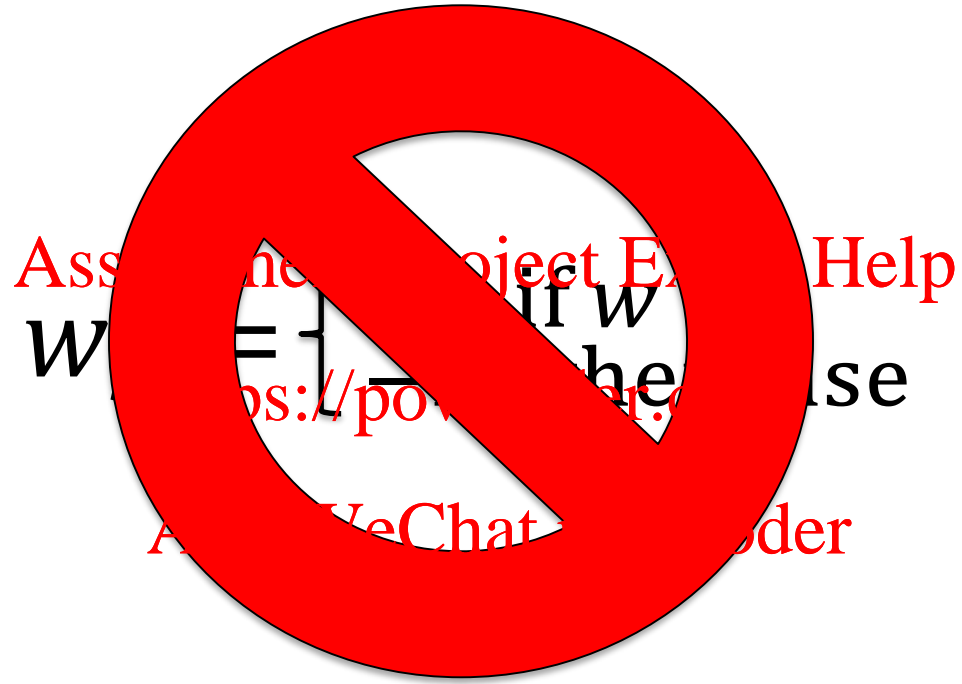
<https://powcoder.com>

Add WeChat powcoder

Runtime comparison



A simple approach



Results in significant loss of information!

Example of binarization

Simple method:

Assignment Project Exam Help
 $w = [0.1, 0.1, 0.1, 0.1]$

$b = [1, 1, 1, 1]$
<https://powcoder.com>

Add WeChat powcoder

A better approach

$$w_b = \begin{cases} +1 & \text{with probability } p = \sigma(w) \\ -1 & \text{with probability } 1 - p \end{cases}$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$\sigma(x) = \max\left(0, \min\left(1, \frac{x + 1}{2}\right)\right)$$

Example of binarization

Reminder: $\sigma(x) = \max\left(0, \min\left(1, \frac{x+1}{2}\right)\right)$

Simple method:

Assignment Project Exam Help

$$w = [0.1, 0.1, 0.1, 0.1]$$

<https://powcoder.com>

$$= [1, 1, 1, 1]$$

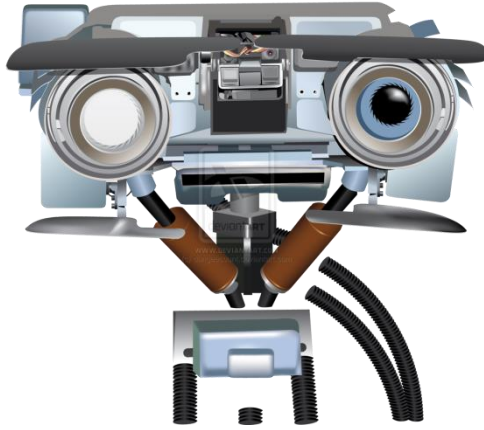
Add WeChat powcoder

Better approach:

$$w = [0.1, 0.1, 0.1, 0.1]$$

$$= [1, -1, -1, 1]$$

Additional details: [Courbariaux et al. BinaryConnect: Training Deep Neural Networks with binary weights during propagations. NeurIPS, 2015.](#)



Assignment Project Exam Help

<https://powcoder.com>

Parameter Pruning

Add WeChat powcoder

CS542 Fall 2020

Remove unimportant weights

- Many weights may not affect performance much (if at all)

Assignment Project Exam Help

- Can save space and computation since you can skip operations*

<https://powcoder.com>

Add WeChat powcoder

Optimal Brain Damage (LeCun et al., NeurIPS, 1989)

- Find parameters that don't affect the training/validation error of the network.
 1. Train your network
 2. Compute the second-order derivatives h_{kk} for each parameter
 3. Compute parameter saliencies $s_k = \frac{h_{kk}u_k}{2}$ ($u_k =$ next layer parameters)
 4. Delete some low saliency parameters
 5. Repeat

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Does this really save space/time?

It depends....

Assignment Project Exam Help

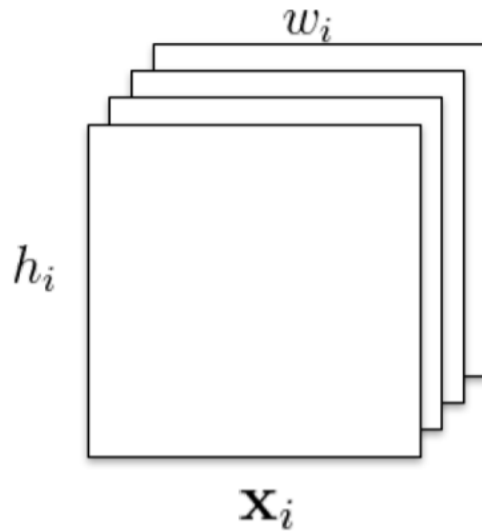
1	1	0
0	1	0
1	1	1

<https://powcoder.com>

Add WeChat powcoder

For CNNs we can just prune entire filters instead!

Pruning entire filters in CNNs



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Pruning process

1. Train your neural network
2. Prune filters for layer l_k
3. Retrain your network
4. Repeat for layer l_{k+1}

Why not all at once?

Assignment Project Exam Help

<https://powcoder.com>

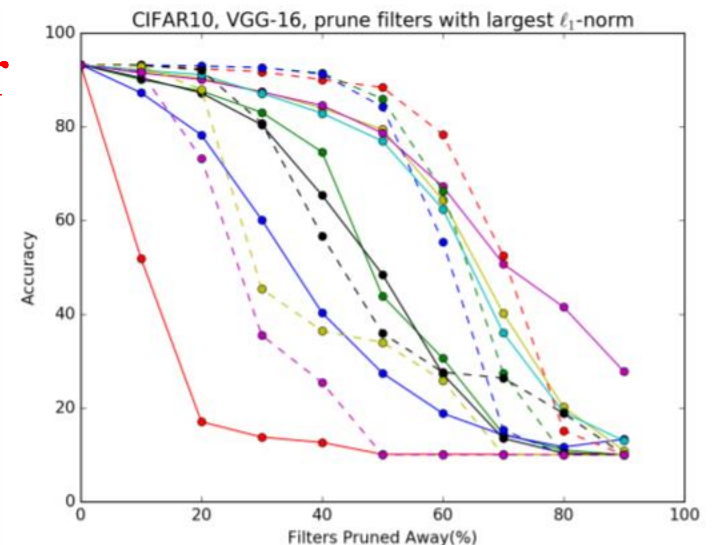
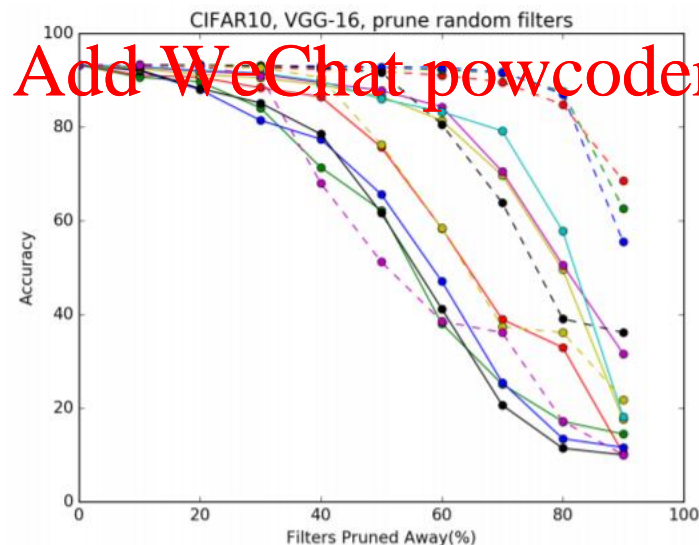
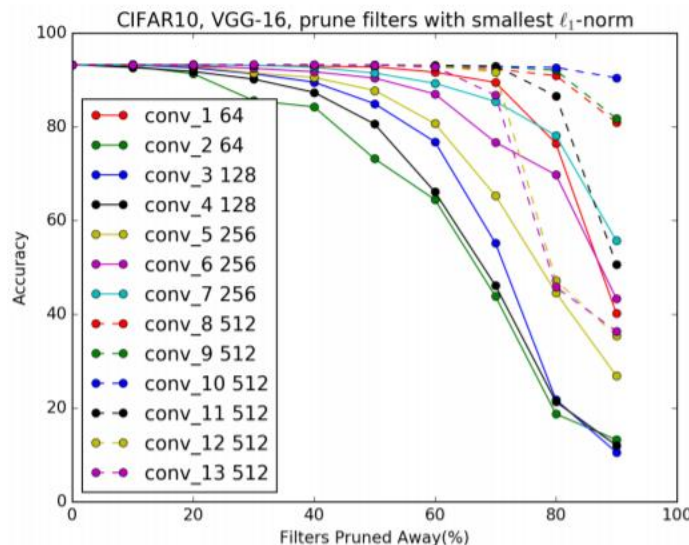
Add WeChat powcoder

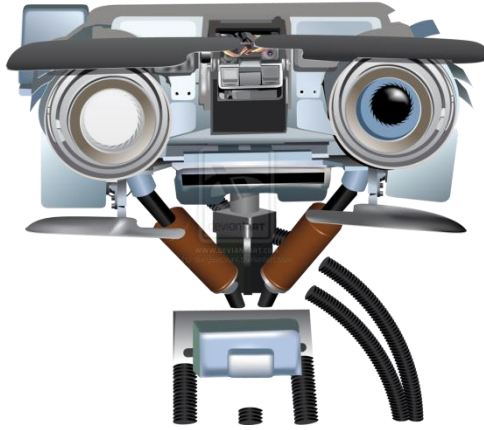
How to choose which filters to prune?

- Often the focus of many research papers.
- The approach of [Li et al., Pruning Filters for Efficient Convnets, ICLR, 2017](#):
– Rank filters w_1, w_2, \dots, w_k using $s_k = \|w_k\|_1$

Assignment Project Exam Help

<https://powcoder.com>





Assignment Project Exam Help

<https://powcoder.com>

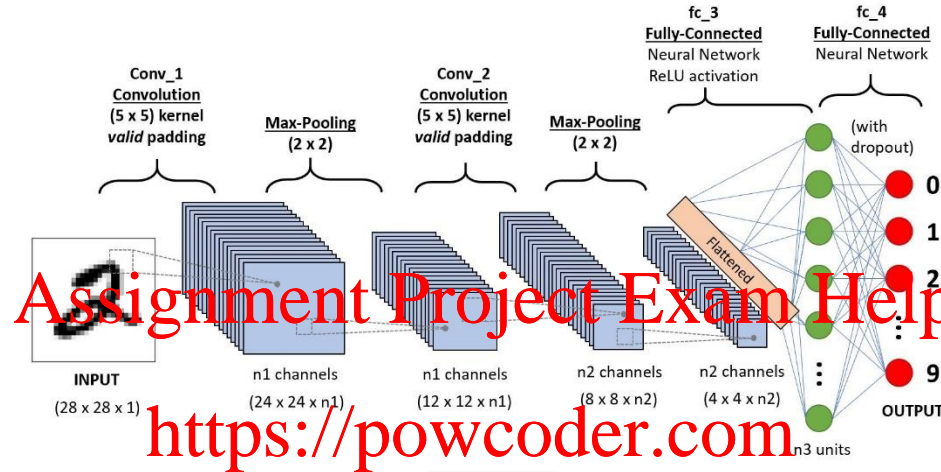
Knowledge Distillation

Add WeChat powcoder

CS542 Fall 2020

Creating an efficient network

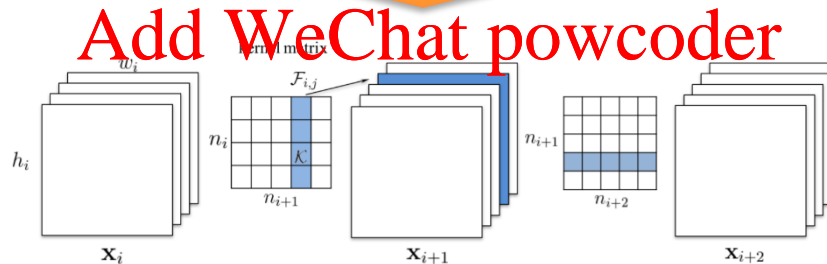
Trained CNN



Assignment Project Exam Help
<https://powcoder.com>

Why not train a more efficient network to begin with?

Pruning

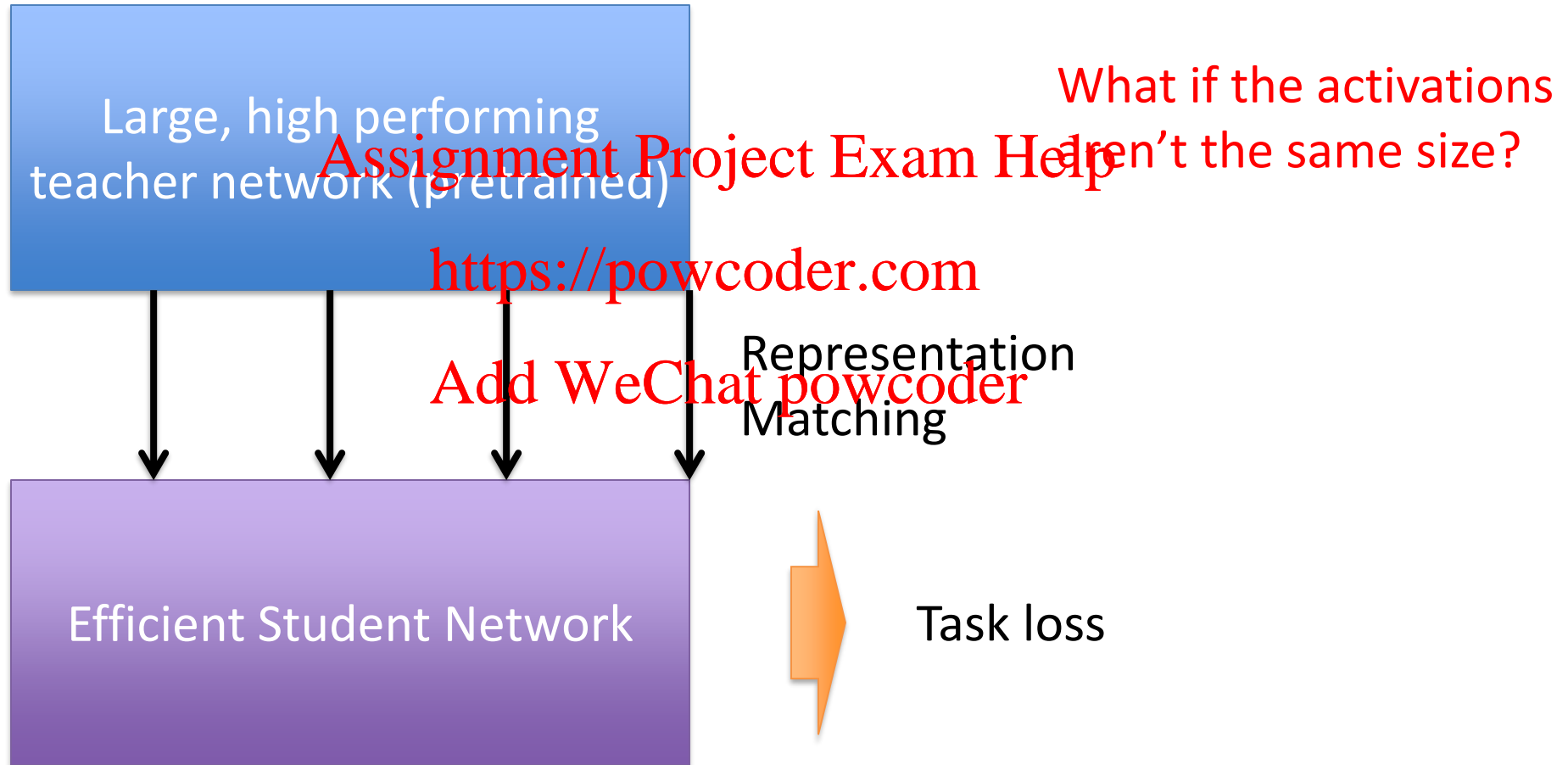


Efficient CNN!

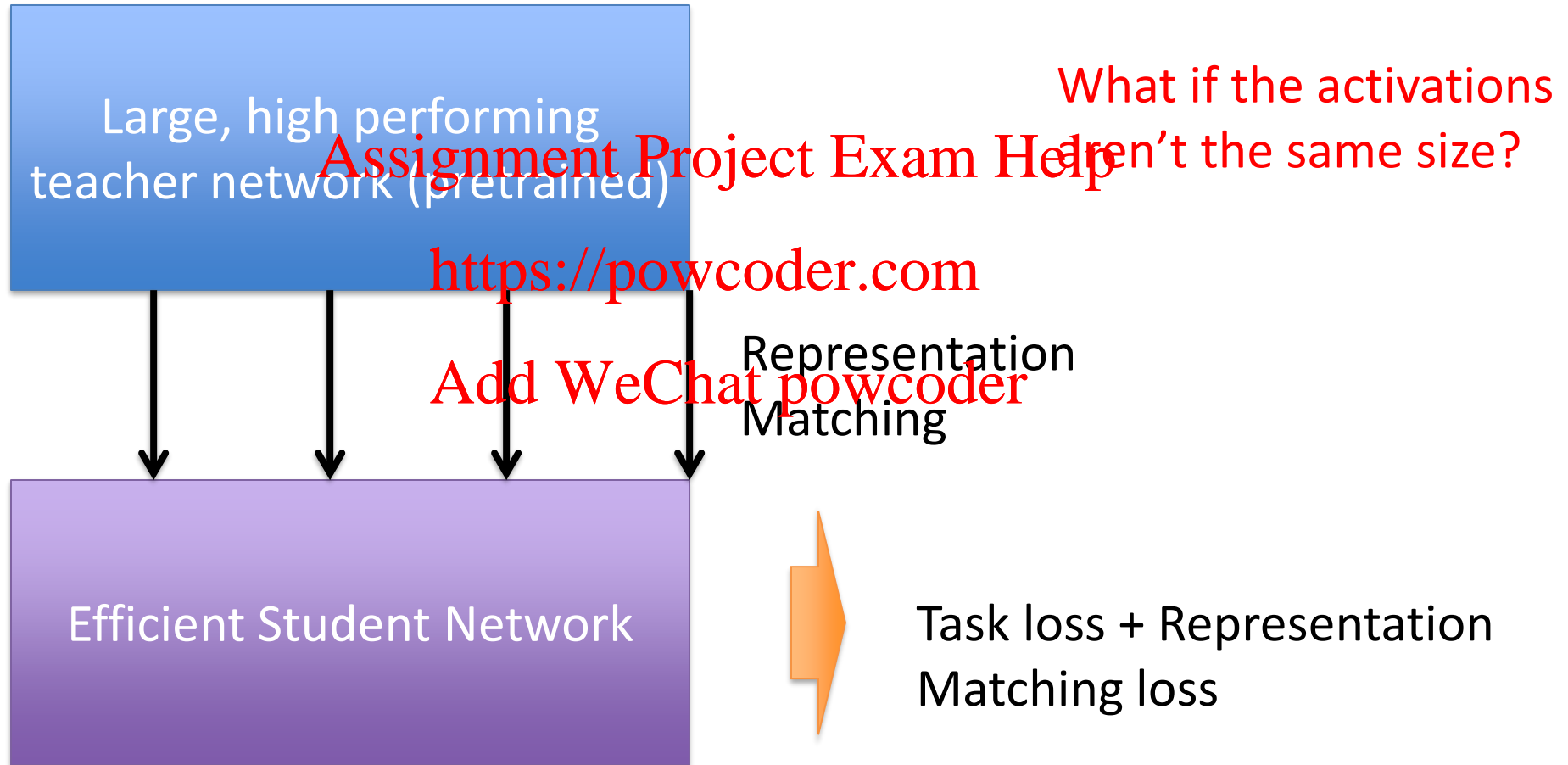
Distillation idea

- Smaller, more efficient networks may not perform well
- We have an example of a network that achieves better performance
<https://powcoder.com>
- Try to get the small network to mimic the large network

Learning with a student-teacher framework



Learning with a student-teacher framework



Distillation loss (i.e. matching representations)

$$L_{distill}(a_t, a_s) = \|a_t - a_s\|_2^2$$

Assignment Project Exam Help

<https://powcoder.com>

More generally considering a case where relevant information is extracted through some transformation function f, g for the teacher, student, respectively

Add WeChat powcoder

$$L_{distill}(a_t, a_s) = \|f(a_t) - g(a_s)\|_2^2$$

Applications of distillation

- Creating an efficient student using a teacher network

Assignment Project Exam Help

- Can distill ensemble results into a single network

<https://powcoder.com>

Method	Speech recognition accuracy
Baseline	58.9
Ensemble of 10 models	61.1
Distilled single model	60.8

Next Class

Final Review:

expect questions on material covered in the entire course in lectures, problem sets, LABs, and assigned reading

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Course evaluations available (5-10 minutes to complete).

Please fill it out at: <https://bu.campuslabs.com/courseeval>