

Network pruning 网络剪枝

因为神经网络中是过参数的, 很多参数没有用.

如何评估一个参数重要与否?

如果绝对值很小, 则说明不重要.

如何评估一个 Neuron 是否重要?

给定 dataset, 如果一个 neuron 的输出基本全为 0, 则说明不重要.

按照重要性排序, 删掉不重要的.

Fine-tune

Q: 为什么不直接 train 小模型, 而用大模型蒸馏?

A: 小模型难 train, 大模型容易优化.

如果裁剪参数 weight 的话, 矩阵运算可能不规则. (将 weight 设为 0)
所以主要裁剪 neuron

知识蒸馏:

Student Net 模仿 Teacher Net 的输出



用 Cross-entropy 来 minimize 输出之间的 loss distribution

Teacher Net 教给学生的更多.

例如: 在识别数字 1 时, 有输出 confidence

1: 0.7 9: 0.1
7: 0.2

- ① 有一个预训练的大网络
- ② 评估重要性
- ③ 删掉不重要
- ④ 微调 Fine-tune
- ⑤ 循环 ②-④ 直到满意

教师学生时会告诉 Student Net 为 1 和 7 长得像。

同时 knowledge distillation 也可以进行 ensemble Learning 的综合。

知识蒸馏的小 trick.

Temperature 方式:

$$y_i = \text{softmax}(X)$$

$$y_i = \text{softmax}(X/T)$$

平滑化, 避免高 confidence
忽略了低 confidence 之间的
相似性

参数量化: ① 用更少的位来^数表征一个值
② 权重聚族.

例如:

1	2	2	1
3	2	1	4
1	4	2	3
4	2	4	3

1, 2, 3, 4 代表 4 种不同范围的
权重, 如 $[0, 1]$ $[2, 3]$

存储位表 将位置处的数值改为
类别编码 (可为 Huffman 编码)
然后存储每种类别的数值 (或其他
方式) 即可

架构设计 =

M 个 neuron

\uparrow weight

N 个 neuron

\Downarrow
共 MN 个 w

M 个

$\uparrow U$

linear K 个 neuron

$\uparrow V$

N 个

\Downarrow
共 $KN + KM$ 个 w
其中 $K < M, N$