

Notes on statistical learning

Adithya Ganesh

December 21, 2018

Contents

1	Matrix cookbook	1
2	Large-scale distributed training	1
3	Hyperparameter optimization	2
3.1	Population-based training	2
3.2	ENAS	2
4	Distillation	2
5	ConvNet architectures	3
5.1	Recognition	3
5.2	Detection	3
5.3	Segmentation	3
5.4	WaveNet	3
5.5	Self-attention networks	3
6	PyTorch	3
7	Backpropagation: CS231 intuitions	3
8	Backpropagation: a graph theory perspective	3
8.1	Combinatorial explosion	5
9	Functional programming \cap Neural networks	6

1 Matrix cookbook

This doesn't matter as much for analytic calculations, but is useful for implementing autodiff on tensors from scratch.

<http://www.math.uwaterloo.ca/~hwolkowi//matrixcookbook.pdf>

2 Large-scale distributed training

Goyal et al. 2017 [?]

- Key contribution: no loss in accuracy when training with large minibatch sizes up to 8192 images.
- Linear scaling rule for adjusting learning rates as a function of mini-batch size. Concretely, when minibatch size is multiplied by k , multiply the learning rate by k .

- Warmup scheme that overcomes optimization challenges early in training. Specifically, use a *low constant* learning rate for the first few epochs of training. For a large minibatch of size kn , train with the low learning rate of η for the first 5 epochs and then return to the target learning rate of $\hat{\eta} = k\eta$.
- Update rule:

$$\hat{w}_{t+1} = w_t - \hat{\eta} \frac{1}{kn} \sum_{j < k} \sum_{x \in \mathcal{B}_j} \nabla l(x, w_t).$$

3 Hyperparameter optimization

3.1 Population-based training

Jaderberg et al. 2017 [?]

- Key contribution: asynchronous optimization algorithm which uses a fixed computational budget to jointly optimize a population of models and their hyperparameters
- Schedule of hyperparameters (instead of fixing a set for the whole course of training)
- Sequential optimization: run multiple training runs (potentially with early stopping)
- Parallel random/grid search: train multiple models in parallel with different weight initializations + hyperparameters, with the view that one of the models will be optimized the best.
- Population based training: starts like parallel search, randomly sampling hyperparameters and weight initializations. Each training run asynchronously evaluates its performance periodically. Explores new hyperparameters by modifying the better model’s hyperparameters, before training is continued.

3.2 ENAS

Pham et al. 2018 [?]

- Key contribution: fast + inexpensive approach for automatic model design.
- Controller (trained with policy gradient) discovers neural network architectures by searching for an optimal subgraph within a large computational graph.
- To train the shared parameters ω of the child models: we fix controller policy $\pi(\mathbf{m}; \theta)$ and perform SGD on ω to minimize expected loss $\mathbf{m} \sim \pi[\mathcal{L}(\mathbf{m}; \omega)]$. Gradient is computed using Monte Carlo estimate

$$\nabla_{\omega} \mathbf{m} \sim \pi(\mathbf{m}; \theta) [\mathcal{L}(\mathbf{m}; \omega)] \approx \frac{1}{M} \sum_{i=1}^M \nabla_{\omega} \mathcal{L}(\mathbf{m}_i, \omega),$$

where \mathbf{m}_i are sampled from $\pi(\mathbf{m}; \theta)$.

- To train controller parameters θ : fix ω and update θ to maximize the expected reward $\mathbf{m} \sim \pi(\mathbf{m}; \theta) [\mathcal{R}(\mathbf{m}, \omega)]$. Use Adam optimizer, and compute the gradient with REINFORCE, with a moving average baseline to reduce variance.

4 Distillation

Hinton et al. 2015 [?]

- Ensembling: train many different models on the same data and then average their prediction
- Distillation: compress the knowledge in an ensemble into a single model which is much easier to deploy.

Anil et al. 2018 [?]

- Online distillation enables extra parallelism.

- Codistillation algorithm:

5 ConvNet architectures

5.1 Recognition

- PNASNet-5-Large
 - Similar to NAS, but performs search progressively (starting with models of low complexity).
- NASNet-A-Large
 - Uses a 50-step RNN as a controller to generate cell specifications.
- SENet154
- PolyNet

5.2 Detection

- Faster RCNN
- YOLO
- RetinaNet

5.3 Segmentation

- FCNet
- DeepLabv4
- Dilated convolutions

5.4 WaveNet

Uses dilated convolutions:

- Let F be a discrete function, and k be a discrete filter. The discrete convolution operator $*$ is defined as

$$(F * k)(\mathbf{p}) = \sum_{\mathbf{s} + \mathbf{t} = \mathbf{p}} F(\mathbf{s})k(\mathbf{t}).$$

More generally, let l be a dilation factor. The l -dilated convolution can be defined as

$$(F *_l k)(\mathbf{p}) = \sum_{\mathbf{s} + l\mathbf{t} = \mathbf{p}} F(\mathbf{s})k(\mathbf{t}).$$

- Implemented in TensorFlow as `tf.nn.atrous_conv2d`

5.5 Self-attention networks

6 PyTorch

7 Backpropagation: CS231 intuitions

8 Backpropagation: a graph theory perspective

Notes from Chris Olah's post: <http://colah.github.io/posts/2015-08-Backprop/>

Algorithm 1 Codistillation

Input loss function $\phi(\text{label}, \text{prediction})$
Input distillation loss function $\psi(\text{aggregated_label}, \text{prediction})$
Input prediction function $F(\theta, \text{input})$
Input learning rate η
for n_burn_in steps **do**
 for θ_i in model_set **do**
 $y, f = \text{get_train_example}()$
 $\theta_i = \theta_i - \eta \nabla_{\theta_i} \{\phi(y, F(\theta_i, f))\}$
 end for
end for
while not converged **do**
 for θ_i in model_set **do**
 $y, f = \text{get_train_example}()$
 $\theta_i = \theta_i - \eta \nabla_{\theta_i} \{\phi(y, F(\theta_i, f)) + \psi(\{\frac{1}{N-1} \sum_{j \neq i} F(\theta_j, f)\}, F(\theta_i, f))\}$
 end for
end while

Algorithm 1 Codistillation

Input loss function $\phi(\text{label}, \text{prediction})$
Input distillation loss function $\psi(\text{aggregated_label}, \text{prediction})$
Input prediction function $F(\theta, \text{input})$
Input learning rate η
for n_burn_in steps **do**
 for θ_i in model_set **do**
 $y, f = \text{get_train_example}()$
 $\theta_i = \theta_i - \eta \nabla_{\theta_i} \{\phi(y, F(\theta_i, f))\}$
 end for
end for
while not converged **do**
 for θ_i in model_set **do**
 $y, f = \text{get_train_example}()$
 $\theta_i = \theta_i - \eta \nabla_{\theta_i} \{\phi(y, F(\theta_i, f)) + \psi(\{\frac{1}{N-1} \sum_{j \neq i} F(\theta_j, f)\}, F(\theta_i, f))\}$
 end for
end while

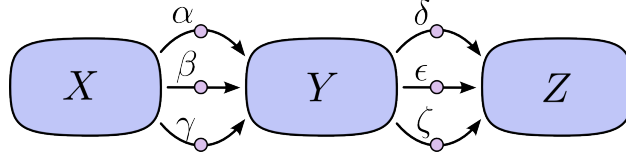
Backpropagation is a very common algorithm, and is often referred to as “reverse-mode differentiation.” At its core, it is a tool for calculating derivatives quickly.

Why are computational graphs a good abstraction? To apply the multivariate chain rule:

1. Sum over all possible paths from one node to the other.
2. Multiply the derivatives on each edge of the path together.

8.1 Combinatorial explosion

This is all just standard chain rule. But how do you deal with cases like this?¹



There are 9 paths in the above diagram. Instead of naively summing over the paths, we can factor them:

$$\frac{\partial Z}{\partial X} = (\alpha + \beta + \gamma)(\delta + \epsilon + \zeta).$$

There are two algorithms we can leverage here.

1. **Forward-mode differentiation.** Start at an input to the graph, and move towards the end. Sum all the paths feeding in. The operator here is $\frac{\partial}{\partial X}$; similar to standard calculus.

$$\begin{aligned}\frac{\partial X}{\partial X} &= 1 \\ \frac{\partial Y}{\partial X} &= \alpha + \beta + \gamma \\ \frac{\partial Z}{\partial X} &= (\alpha + \beta + \gamma)(\delta + \epsilon + \zeta).\end{aligned}$$

2. **Reverse-mode differentiation.** Start at an output of the graph, and move towards the beginning. At each node, merge all paths which started at that node. The operator here is $\frac{\partial Z}{\partial}$.

In particular:

$$\begin{aligned}\frac{\partial Z}{\partial Z} &= 1 \\ \frac{\partial Z}{\partial Y} &= \delta + \epsilon + \zeta \\ \frac{\partial Z}{\partial X} &= (\alpha + \beta + \gamma)(\delta + \epsilon + \zeta)\end{aligned}$$

What is the difference between forward and reverse mode differentiation? “Forward-mode differentiation tracks how one input affects every node.” “Reverse-mode differentiation tracks how every node affects one output.”

¹Image source: Chris Olah.

$$\begin{array}{l} \text{forward mode: } \frac{\partial}{\partial X} \\ \text{reverse mode: } \frac{\partial Z}{\partial} \end{array}$$

Reverse mode differentiation gives us the derivative of the output w.r.t. every node. This is exactly what we want. On a large computational graph, this means reverse mode differentiation can get them all in one fell swoop. In summary: derivatives are ridiculously computationally cheap.

9 Functional programming \cap Neural networks

TODO: write up notes on Chris' article.