March 2020

# Designing NETWORK Design Spaces

So these guys come in and say
hey up until now we have been
looking at a model space then
finding the best singular model
ether Manually or automaticly
Within that space. But they say
lets take it a step/layer
above And find a space Which
describes a subset of all models
Where these models on avg
Better But also "simpler, work well
and generalize across settings"

They Start With an rather
unconstrained space and progressivly
Constraing it "while maintaing or imporoving"
the error distribution produced by the
Models.

The least Constrained Space is
Called AnyNET and is as follows:
(note look at paper for good drawings)



The Body is where they will
Be defining the model + the Bulk of the
Work Will Be done. Stem = Stide2 3×3 conv
32 channels + the head is AvG Pooling +a FC luyer
There are 4 Stages in the Body Where
each stage $i$ has $b_i$ (blocks) $w_i$ (width) + other
Block Params.

Since each network has 4 stages
& each stage has 4 degrees of
freedom in total there are 16 degrees
of freedom.

each stage has $d_i$ (Blocks)
$w_i$ (width), $b_i$ (Bottleneck ratio), and $g_i$ (group width)

$$d_i \leq 16 \qquad w_i \leq 1024 \text{ (and divisable by 8)}$$

$$b_i = \{1, 2, 4\} \qquad g_i \in \{1, 2, \ldots 32\}$$

So above is the AnyNet design space
with $10^{18}$ possible model configs.

Step one they set all $b_i = b$
So it's the same across all stages of
a model. They find no increas in
error But now the design space
is simpler.

Step two is to set $g_i = g$
as above & they find same result

Step 3 they find Pattern after
Step 2 where increasing width over the
stages results in Better models
So they test AnyNet$_D$ where AnyNet$_C$
is after Step 2 & AnyNet$_D$ is
after Step2 & only models where
$w_{i+1} \geq w_i$ & find it significantly
Better distribution of error

Step 4    they    find    that    similar
as   with   step 3   if   now   we
increase   depth   $d_{i+1} \geq d_i$   the   models
are   better.

So   after   all   these   reductions
our   design   space   went   from
$10^{18}$ posobility   to   $10^7$

So   then   they   come   up   with   the
final   Design space   Described   as   such:

RegNET   generated   from:   $d, W_0, W_a, W_m$
$d < 64$    $W_0, W_a \leq 256$   But we Have
                         to  qhantize
$1.5 \leq W_m \leq 3$    $U_j = W_0 + W_a \cdot j$   via   Control the
                                             $W_m$   Scaling of
                                                     width

The   original   tests   we   have   Been
reading   about   are   all   done   in   the
Low   epochs   +   low   compute range.
So   now   they   compare   in   higher   compute
higher   epoch   +   5 stages
the   ordering   is   always   $RegNetx > AnyNet_E > AnyNet_A$
                                                    Post      Pre
                                                   Step 4    Step 1

They   then   have   further   observations
that   the   common   $b < 1$   +   $g = 1$
are   not   as   good   as   $b = 1$   +   $g \geq 1$

also

they also found optimal depth = 20
Blocks (intwessing Deeper not always Better!)

& a width multipu of 2.5 close to the
common one of 2

So Now lets compare ResNET
model to other models.

RegNET Models tend to have
lower Flops But maintain or
Better sresults the ResNet

In general the RegNETS
Matched or did Better than
state of the art ResNET

and at low flop Efficent net
Better but at Higher flops RegNET
Better. & is much faster in the
Higher flop regions

$d = 4$

Say: $W_0 = 32$

$W_a = 8$

$W_m = 2$

On quantization: So we get Powers of 2

| | | | |
|---|---|---|---|
| $32 = U_0 = 32 + 8 \cdot 0$ | $\log\left(\frac{32}{32}\right) = 0$ | Round $= 0$ | $W_0 = 32 \cdot 2^0 = 32$ |
| $46 = U_1 = 32 + 8 \cdot 1$ | $\log(40/32) = .32$ | $\Rightarrow = 0$ | $W_1 = 32 \cdot 2^0 = 32$ |
| $48 = U_2 = 32 + 8 \cdot 2$ | $\log(48/32) = .58$ | $= 1$ | $W_2 = 32 \cdot 2^1 = 64$ |
| $56 = U_3 = 32 + 8 \cdot 3$ | $\log(56/32) = .81$ | $= 1$ | $W_3 = 32 \cdot 2^1 = 64$ |

Now 2 Stages

|  Stage 1 | Stage 2 |
|---|---|
| 2 Blocks | 2 Blocks |
| 32 width | 64 width |