

Regularisation Methods

[Neural Networks]

- Adam
- L1 / L2
- Dropout
- Batch Norm

Regulations

in L_R , we use

$$L: \min_{\vec{w}, \gamma} \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2 + \lambda (\|\vec{w}\|_2)$$

$$\rightarrow L: \max_{\vec{w}, \gamma} f(\vec{w}) + \underbrace{\lambda (\|\vec{w}\|_2^2)}_{\text{Regularization}}$$

We can do something similar but here W is a matrix.

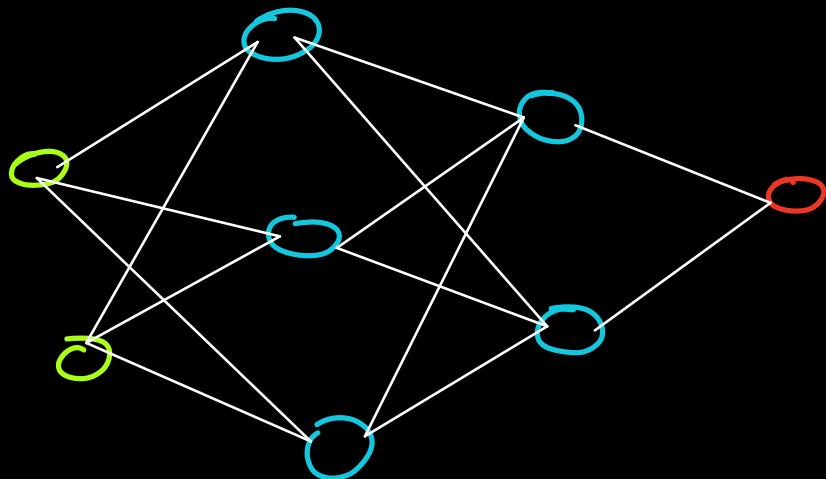
$$\text{Reg : } \lambda \sum_{i=1}^m \sum_{j=1}^n |w_{ij}|^2 \quad (\text{for all layers})$$

Frobenius Norm

(Refer to notes for grand)

```
def create_baseline():
    # lambda = 0.01
    L2Reg = tf.keras.regularizers.L2(l2=1e-6)
    model = Sequential([
        Dense(256, activation="relu", kernel_regularizer = L2Reg ),
        Dense(128, activation="relu", kernel_regularizer = L2Reg),
        Dense(64, activation="relu", kernel_regularizer = L2Reg),
        Dense(1 , activation = 'sigmoid')])
    return model
```

Drop Out



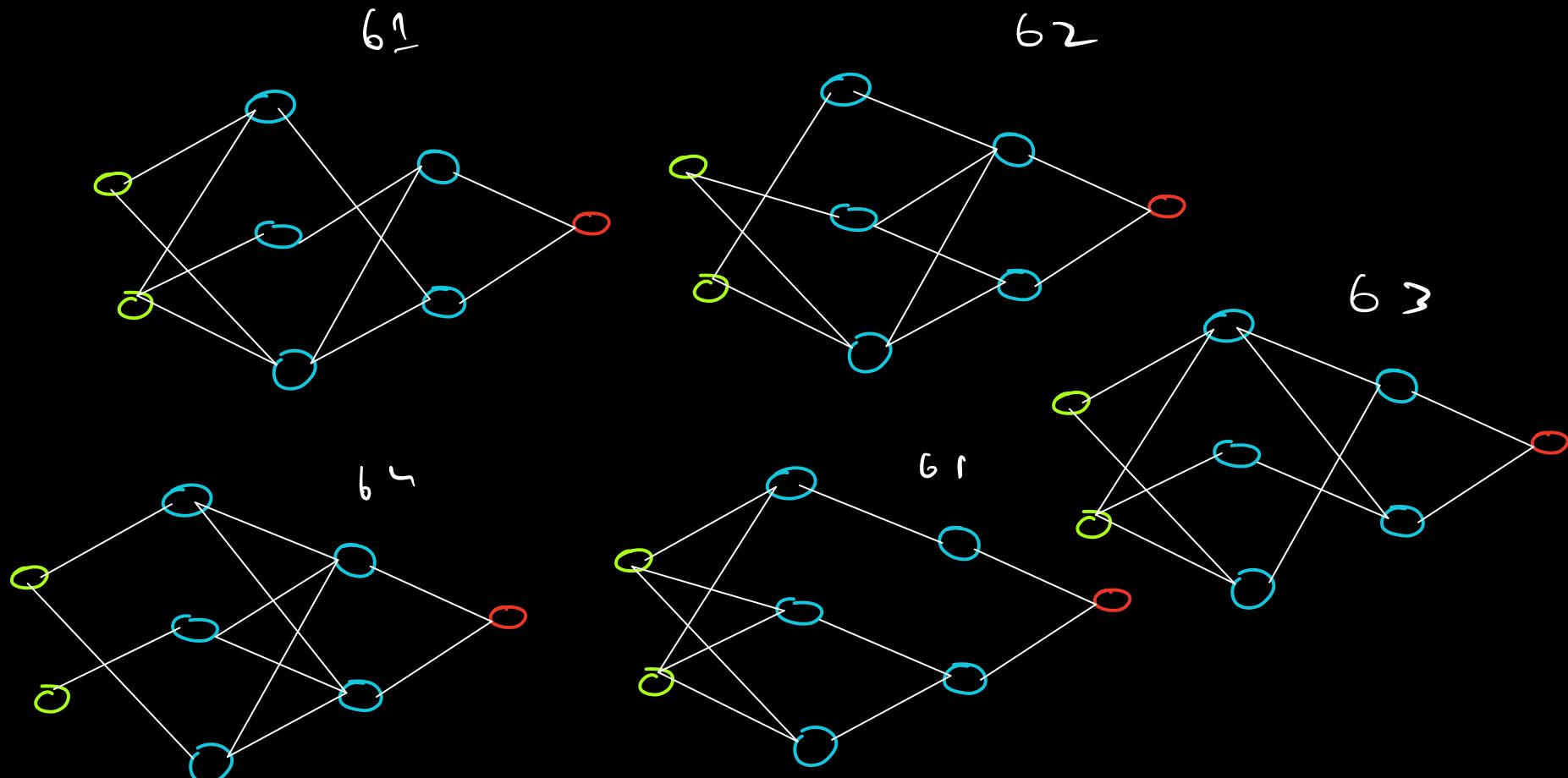
$$w_s \rightarrow 2 \times 3 \\ b_s \rightarrow +3$$

$$3 \times 2 \\ +2$$

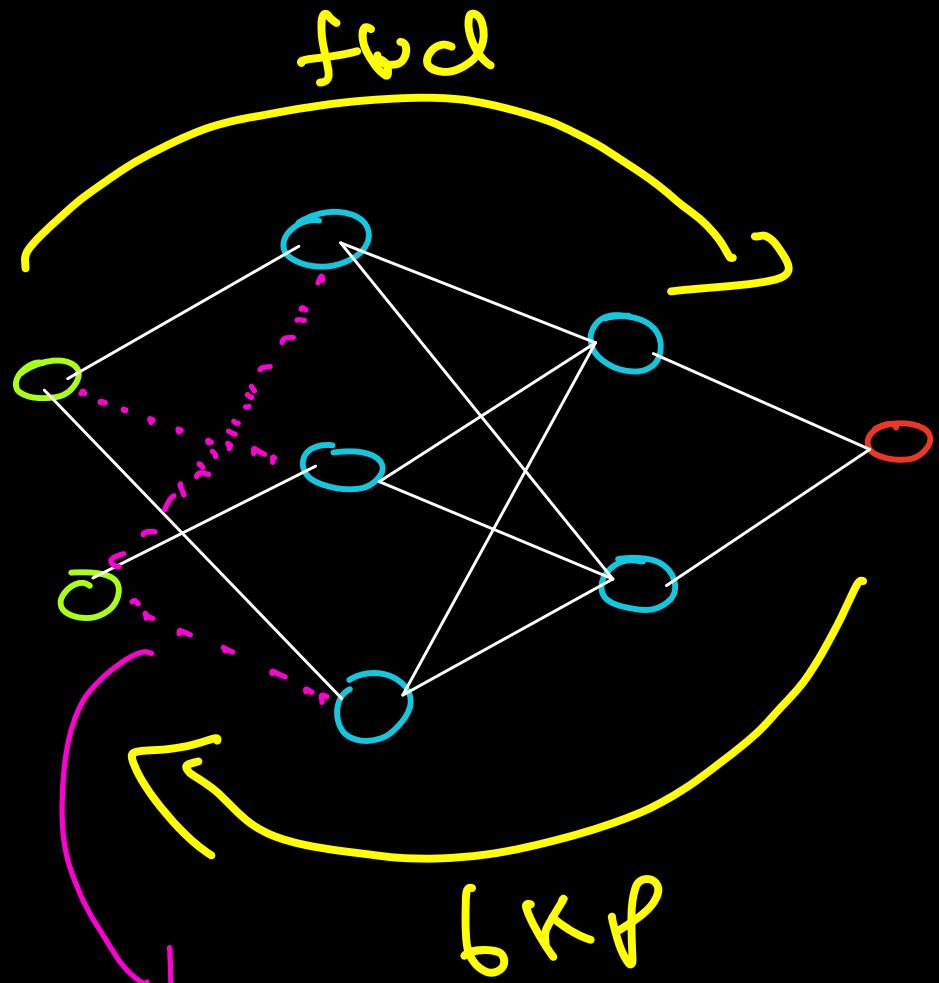
$$\begin{array}{rcl} 2 & = 14 \\ +1 & = 6 \\ \hline 20 & & \text{persons} \checkmark \end{array}$$

With deep or NN # persons becomes very high. It is common to have more persons than training points. This causes overfitting.

Solution



→ Randomly
during each batch.
remove (drop) connections



missing weights won't be
updated in this iteration

This "kind of"
creates an ensemble
of many "partially"
connected NN.

We know that
bagging reduces
overfitting (variance)

Each update has less weights now.
so less overfits.

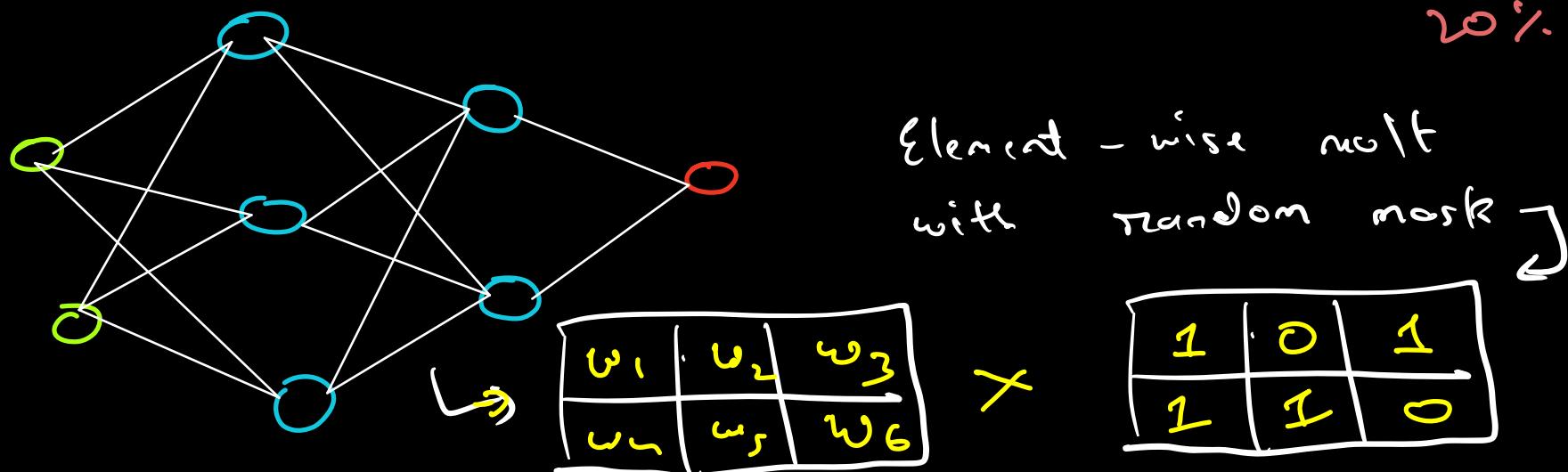
→ Plus the model needs to perform well for diff combination of weights

How is this implemented?

HP → specify dropout rate, eg: 0.2

↓

20%



Dropout layer has "no trainable"
parameters

During Inference

Since training is done with only
e.g. 80% of the weights at any
given time,

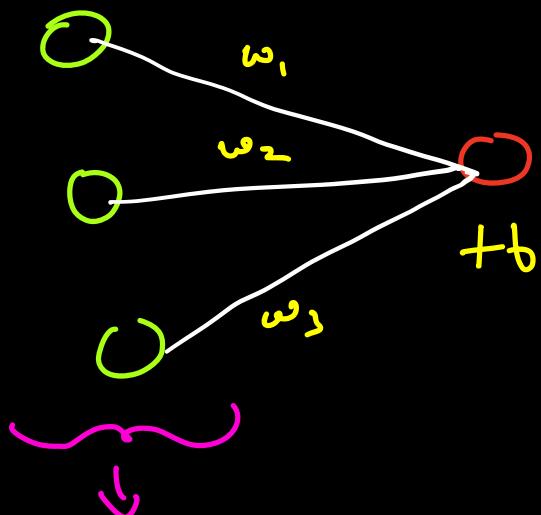
During prediction we use 100% of
the weights, but we multiply each
weight by say 0.8

Dropout does not affect biases.

```
from tensorflow.keras.layers import Dropout
def create_Dropout():
    # lambda = 0.01
    L2Reg = tf.keras.regularizers.L2(l2=1e-6)
    model = Sequential([
        Dense(256, activation="relu", kernel_regularizer = L2Reg ),
        Dropout(0.3),
        Dense(128, activation="relu", kernel_regularizer = L2Reg),
        Dropout(0.3),
        Dense(64, activation="relu", kernel_regularizer = L2Reg),
        Dense(1 , activation = 'sigmoid')])
    return model
```

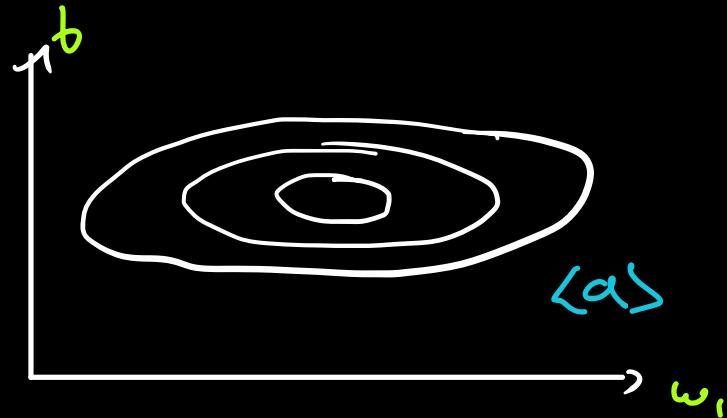
Batch Normalisation

Let's consider a simple NN

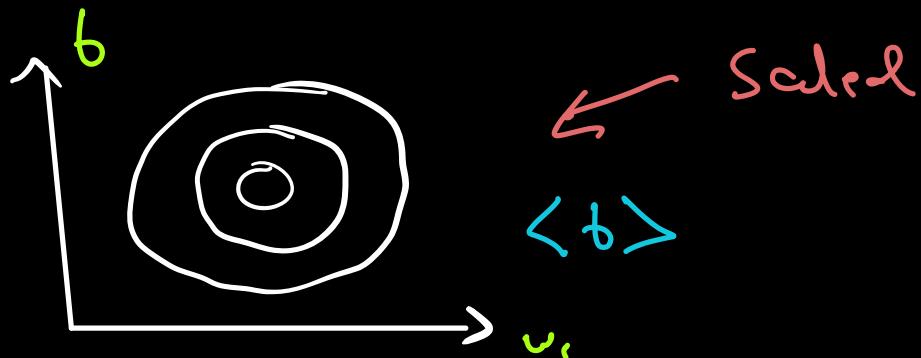


Should I
scale these
values?

$$\# \text{ params} = 4$$



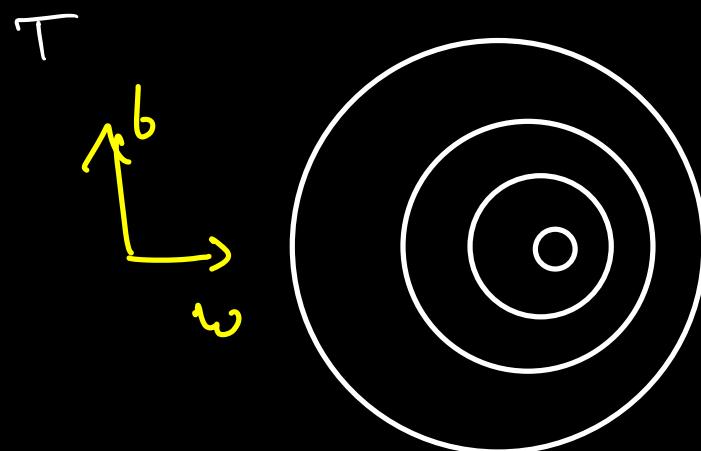
Q. Where
will it
train
faster?



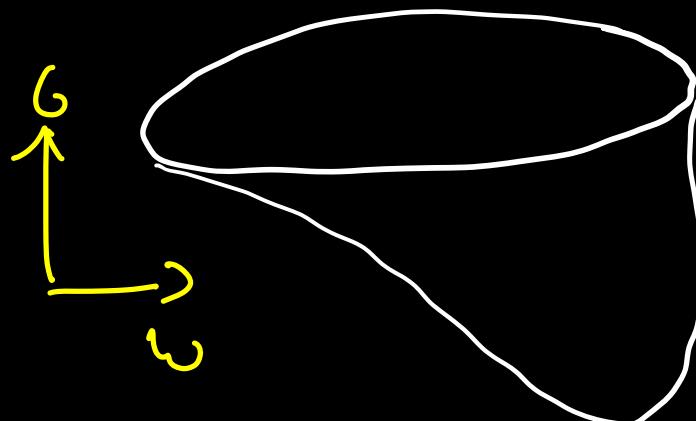
- if we scale the input, that makes w and b ranges similar to each other.
- It is well known that standardisation makes training faster!

But,

You still need RMSProp / ADAM



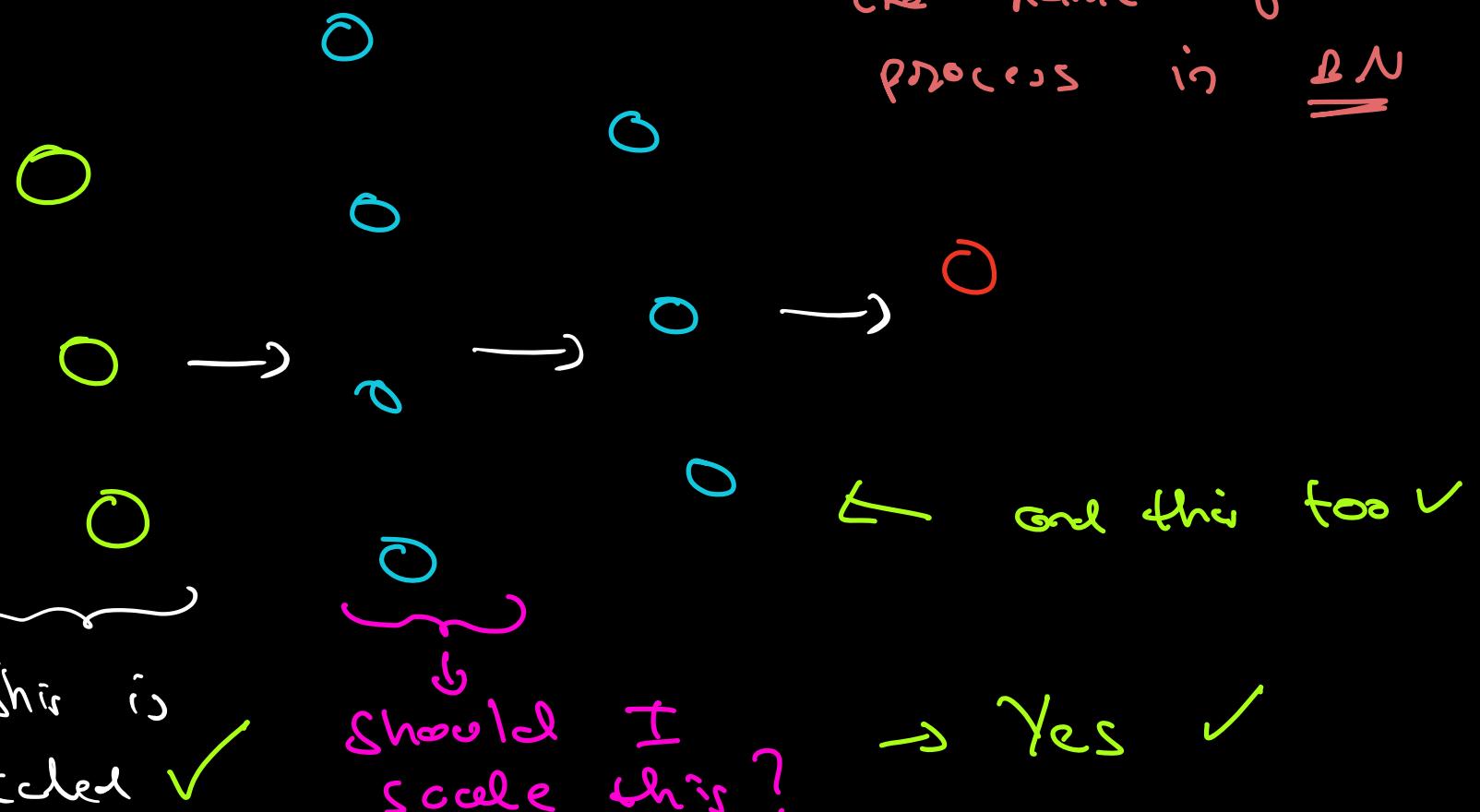
Contours [Top View]



d_w is very large on this side

So a combintⁿ of standardisation and
optimised helps!

Further



Std scaling the output of a layer is
over every batch is called Batch Normalization

There is still debate on whether to apply
it before or after the activation func

$$z' \xrightarrow{BN} \tilde{z}' \xrightarrow{a'} \langle a \rangle$$

OR

$$z' \xrightarrow{a'} \tilde{a}' \xrightarrow{BN} \langle b \rangle$$

→ You can try this out via code ses.

```

def create_BatchNormalization_model():
    L2Reg = tf.keras.regularizers.L2(l2=1e-6)
    model = Sequential([
        Dense(256, kernel_regularizer = L2Reg, activation='relu'),
        BatchNormalization(), # <b>
        Dense(128, kernel_regularizer = L2Reg, activation=None),
        BatchNormalization(), # <a>
        Activation('relu'),
        Dense(64,kernel_regularizer = L2Reg ),
        BatchNormalization(),
        Activation('relu'),
        Dense(1)])
    return model

```

How does this work?

1 Batch with 3 samples mean std-dev					
Features	x_1	x_2	x_3	x_4	
	1	3	8	4	2.94
	3	4	3	3.33	0.471
	5	6	2	4.33	1.69
	7	2	1	3.33	2.62

Normalization across mini-batch,
independently for each feature

We scale one batch at a time!
which means diff batches may have diff μ , σ .

$$m_b = \frac{1}{B} \sum_{i=1}^B x_i \quad ; \quad B = \text{batch size}$$

$$\hat{\sigma}_b^2 = \frac{1}{B} \sum_{i=1}^B (x_i - m_b)^2$$

$$\hat{x}_i = \frac{x_i - m_b}{\sqrt{\hat{\sigma}_b^2 + \epsilon}} \rightarrow \text{to protect from div by 0}$$

$$y_i = BN(x_i) = \gamma \cdot \hat{x}_i + \beta$$

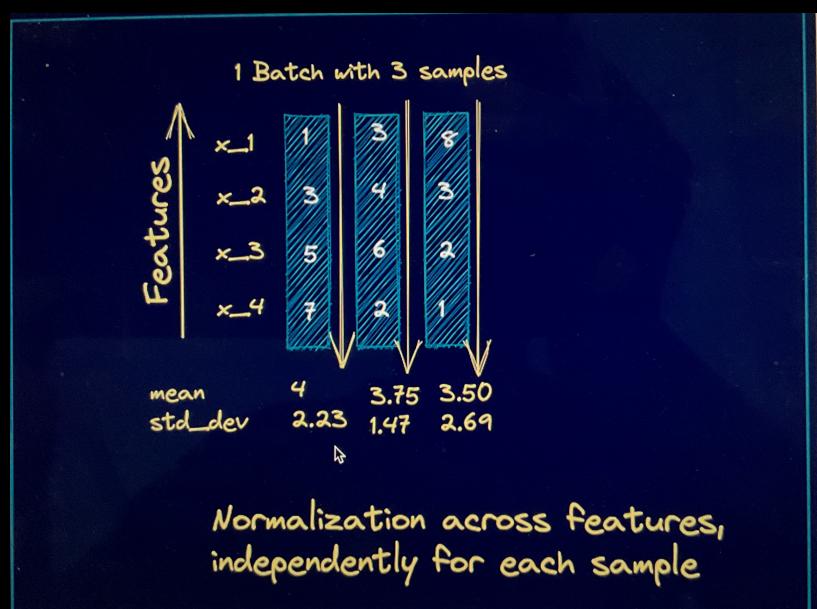
γ, β
additional params!

If each layer has the output with 0 mean and 1 σ then, we may

not learn much!

So we allow the model to figure out "automatically" the desired mean (β) and std (δ)

Please note that there is also something called Layer Normalization



But how does this help in regularization?

→ Each batch has a different μ_b & σ_b . This forces the models to learn such weights that perform well on each batch (diff μ, σ) equally (on avg), hence forcing it to generalize better.

Finally,
what μ, σ to use during inference??

→ we take an exponentially weighted avg of all the μ, σ learned over all batches over all epochs.