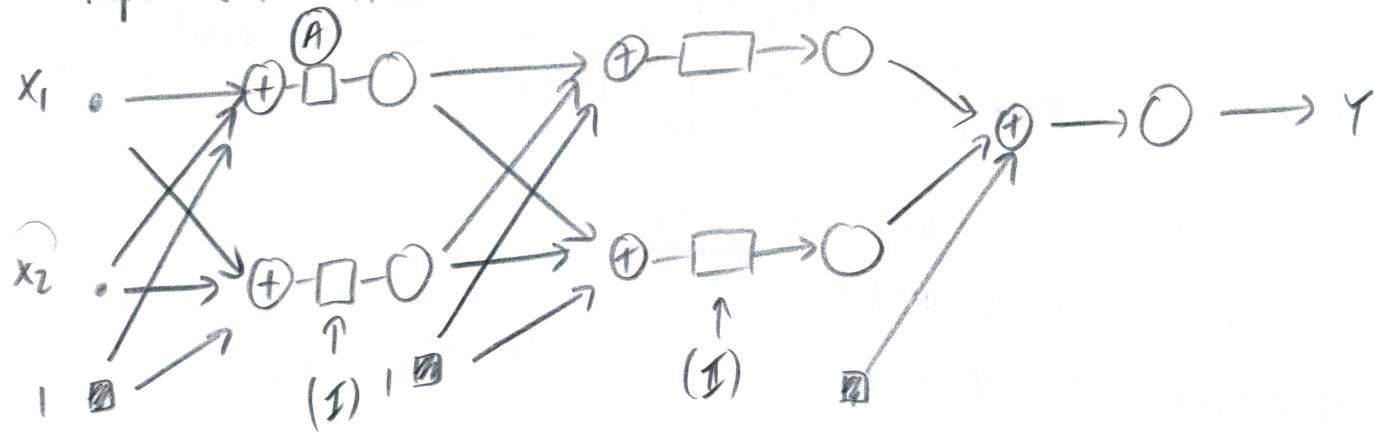


L7 - SGD, Overfitting, regularisation (BN part)- Batch Normalisation (part of HW1)- diagram- 2 hidden layer MLP

- BN - covariate adjust. unit that happens after the weighted addition of inputs (after affine transformations); and before the activations.



Note BN related transforms at (I) and (II).

Individual/close-up view

- BN happens at an individual unit

- Training → adjustment occurs over individual minibatches

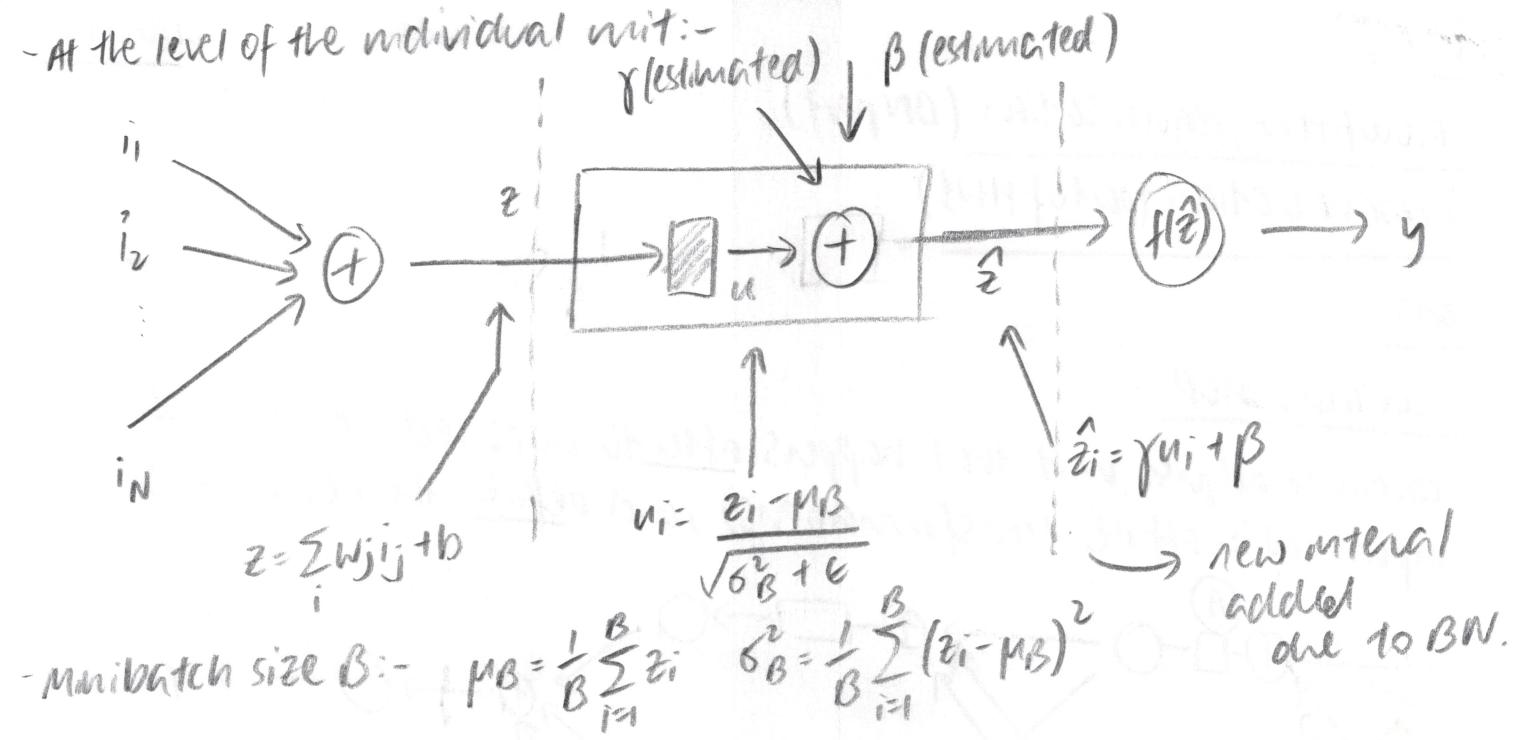
- BN aggregates stats over a minibatch at a specific neuron, normalises

1) every instance within the mini-batch at the particular unit  
will end up being normalised by statistics computed over the  
mini-batch

2) A further linear transformation is applied on every instance  
normalised instance within the mini-batch.

(to a unit-specific location)

② - see supplementary notes on Batch Norm; and Ioffe and Szegedy (2015)



(\*) next piece of notation in lectures is confusing

(\*) don't B to refresh memory that:-

(3) i) Batch gradient descent:-

$$\underline{w}_R \leftarrow \underline{w}_R - \eta \left( \frac{1}{N} \sum_{t=1}^N \nabla_{\underline{w}_R} \text{div}(y_t, d_t) \right)$$

ii) Incremental update (with SGD)

$$\underline{w}_R \leftarrow \underline{w}_R - \eta \left( \nabla_{\underline{w}_R} \text{div}(y_t, d_t) \right)$$

iii) Mini-batch gd.

$$\underline{w}_R \leftarrow \underline{w}_R - \eta \left( \frac{1}{m} \sum_{t=1}^m \nabla_{\underline{w}_R} \text{div}(y_t, d_t) \right)$$

$B/m$ -size of mini-batch

(\*) Q/S1 - note mini-batches is analogous to statistical sampling

↳ more details when clarified

(ii) more batchnormalisation; the function  $\text{div}(y_t, d_t)$  takes additional arguments  $\mu_B$  and  $\sigma_B^2$

i.e. we have  $\left( \frac{1}{m} \sum_{t=1}^m \nabla_{\underline{w}_R} \text{div}(y_t(x_t, \mu_B, \sigma_B^2), d_t) \right)$

(\*) lecture notes state that  $d_t = d_t(x_t)$ ; but I'm not sure now the target label  $d_t$  can be a function of input.

Typo?

## @key point

- $\mu_B$  and  $\sigma_B^2$  are functions of all the other training instances within the mini-batch (they are statistics)
- i.e.  $\mu_B = f(x_1, \dots, x_m)$  and  $\sigma_B^2 = g(x_1, \dots, x_m, \mu_B)$
- Rendering:-  

$$\frac{1}{m} \sum_{t=1}^m J_{WR} \text{div}\left(\gamma_t(x_t, \mu_B(x_1, \dots, x_m), \sigma_B^2(x_1, \dots, x_m, \mu_B)), d_t\right)$$

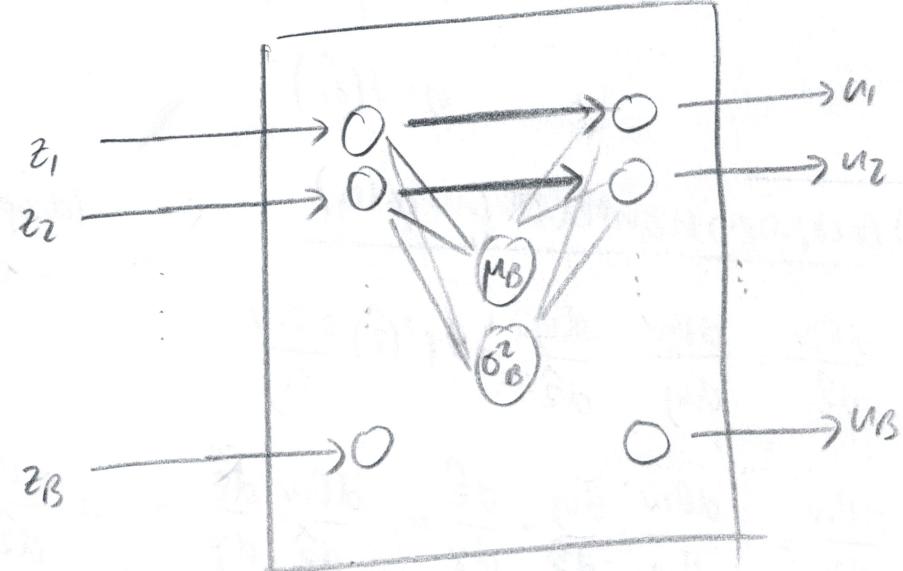
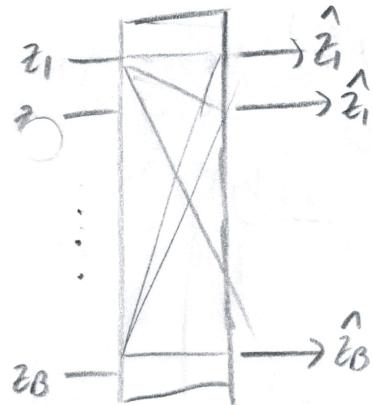
(\*) This complicates backpropagation; but not to the degree of intractability

(\*) Batchnorm as a vector function over a minibatch

(\*) Note it is not a vector function over dimensions (as we saw with softmax)

- It is a vector function over all inputs from a mini-batch
- computing  $\frac{\partial \text{div}}{\partial z_i}$   $\Rightarrow$  consider all  $\hat{z}_j$ 's

solve influence diagrams:-

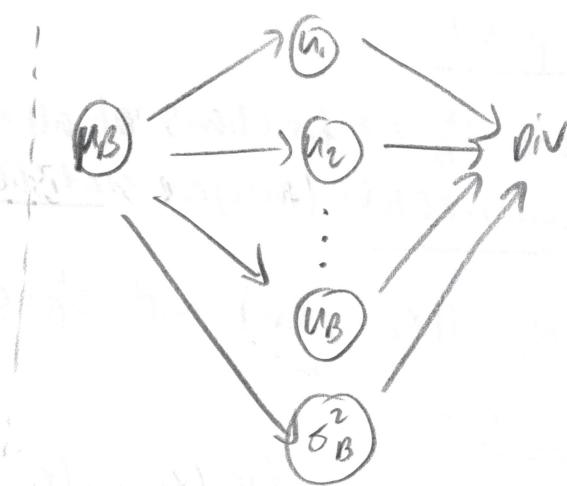
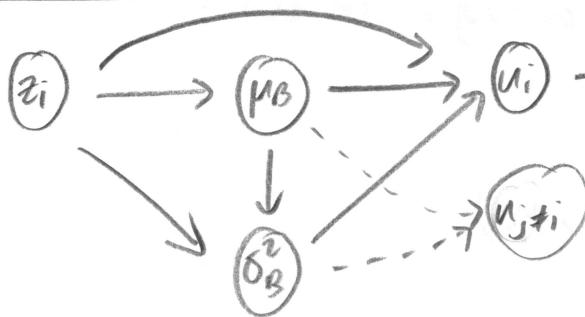


(\*) Inputs/outputs are different instances in a minibatch

(\*) BN occurring at level of individual neuron

(\*) Use these to aid understanding; vector function deriv to assist with next diagram

## Influence diagrams



## (\*) BN-backpropagation

- start at end; move backwards recursively; formed first  
 (for one instance within batch) add at n div. neuron

### NO BN

$$z_i = \sum_j w_{ij} j + b$$

$$y_i = f(z_i)$$

### with BN

$$z_i = \sum_j w_{ij} j + b$$

$$u_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\mu_B = \frac{1}{B} \sum_{i=1}^B z_i \quad \sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (z_i - \mu_B)^2$$

$$\hat{z}_i = \gamma u_i + \beta$$

$$y = f(\hat{z}_i)$$

(drop all subscripts)

- to preserve  
emphasis  
on generality

## Backprop → high level (no calc.)

$$\frac{d\text{Div}}{d\hat{z}} = \frac{d\text{Div}}{dy} \cdot \frac{dy}{d\hat{z}} = f'(\hat{z}) \frac{d\text{Div}}{dy}$$

$$\frac{d\text{Div}}{dy} = \underbrace{\frac{d\text{Div}}{dy} \cdot \frac{dy}{d\hat{z}} \cdot \frac{d\hat{z}}{dy}}_{\text{full calc.}} = \underbrace{\frac{d\text{Div}}{d\hat{z}}}_{\text{reversible part}} \frac{d\hat{z}}{dy} = u \frac{d\text{Div}}{d\hat{z}}$$

$$\frac{d\text{Div}}{d\beta} = \frac{d\text{Div}}{d\hat{z}} \cdot \underbrace{\frac{d\hat{z}}{d\beta}}_{=1} = \frac{d\text{Div}}{d\hat{z}}$$

$$\frac{d\text{Div}}{du} = \frac{d\text{Div}}{d\hat{z}} \cdot \underbrace{\frac{d\hat{z}}{du}}_{=f'} = f' \frac{d\text{Div}}{d\hat{z}}$$

- we have:-

$$\frac{\partial \text{Div}}{\partial \mu_B} = \sum_{i=1}^B \frac{\partial \text{Div}}{\partial u_i} \frac{\partial u_i}{\partial \mu_B} + \frac{\partial \text{Div}}{\partial \sigma_B^2} \cdot \frac{\partial \sigma_B^2}{\partial \mu_B} ; \quad u_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad (a)$$

$$\sigma_B^2 = \frac{1}{B} \sum_{i=1}^B z_i^2 - 2z_i \mu_B + \mu_B^2 \quad (b)$$

From (a) :-  $\frac{\partial u_i}{\partial \mu_B} = \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}}$

From (b) :-  $\frac{\partial \sigma_B^2}{\partial \mu_B} = \frac{\partial}{\partial \mu_B} \frac{1}{B} \sum_{i=1}^B z_i^2 - 2z_i \mu_B + \mu_B^2$   
 $= \frac{1}{B} \sum_{i=1}^B \frac{\partial}{\partial \mu_B} (z_i^2 - 2z_i \mu_B + \mu_B^2) = \frac{1}{B} \sum_{i=1}^B -2z_i + 2\mu_B$

∴  $\frac{\partial \sigma_B^2}{\partial \mu_B} = \frac{\sum_{i=1}^B -2(z_i - \mu_B)}{B}$

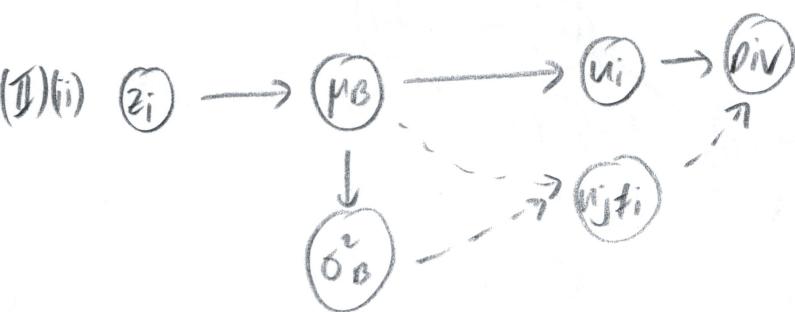
∴  $\frac{\partial \text{Div}}{\partial \mu_B} = \sum_{i=1}^B \frac{\partial \text{Div}}{\partial u_i} \left( \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial \text{Div}}{\partial \sigma_B^2} \left( \frac{\sum_{i=1}^B -2(z_i - \mu_B)}{B} \right)$   
 $= 0$

(more explicitly)

Note:  $\frac{\sum_{i=1}^B -2(z_i - \mu_B)}{B} = -2 \sum_{i=1}^B \frac{z_i - \mu_B}{B} = -2 \left\{ \sum_{i=1}^B \frac{z_i}{B} - \sum_{i=1}^B \frac{\mu_B}{B} \right\} = -2 \{ \mu_B - \bar{z} \}$

(\*)

Hence:  $\frac{\partial \text{Div}}{\partial \mu_B} = \sum_{i=1}^B \frac{\partial \text{Div}}{\partial u_i} \left( \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right)$

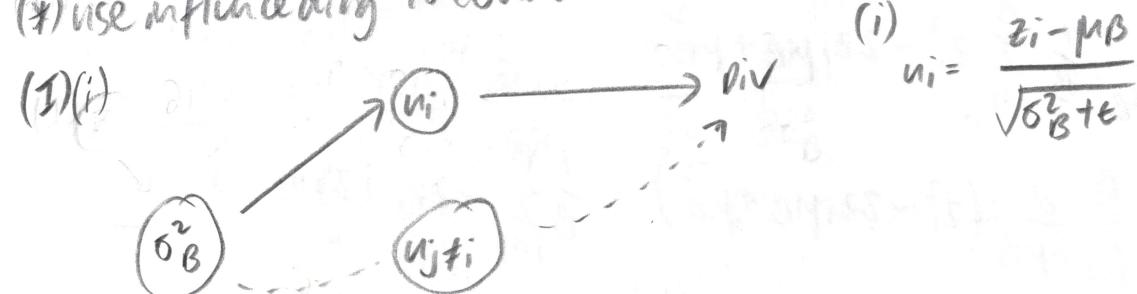


$\frac{\partial \mu_B}{\partial z_i} = \frac{1}{B}$

• And more trickily, but not intractably (removing subscripts)  
 →  $\text{div} = f(u_i, \mu_B, \sigma_B^2)$  →  $\text{div}$  must read this with INFLUENCE DIAG.

$$\frac{\partial \text{div}}{\partial z_i} = \frac{\partial \text{div}}{\partial u_i} \cdot \frac{\partial u_i}{\partial z_i} + \frac{\partial \text{div}}{\partial \sigma_B^2} \cdot \frac{\partial \sigma_B^2}{\partial z_i} + \frac{\partial \text{div}}{\partial \mu_B} \cdot \frac{\partial \mu_B}{\partial z_i}$$

(\*) use influence diag to evaluate each term:-



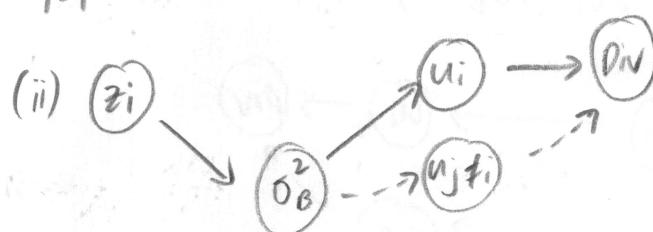
$$\frac{\partial \text{div}}{\partial z_i} = \sum_{j=1}^B \frac{\partial \text{div}}{\partial u_j} \frac{\partial u_j}{\partial \sigma_B^2} + \underbrace{\frac{\partial \text{div}}{\partial u_i} \frac{\partial u_i}{\partial \sigma_B^2}}_{\text{direct}} = \sum_{i=1}^B \frac{\partial \text{div}}{\partial u_i} \frac{\partial u_i}{\partial \sigma_B^2}$$

from (i);  $\frac{\partial u_i}{\partial \sigma_B^2} = -\frac{1}{2}(z_i - \mu_B)(\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$

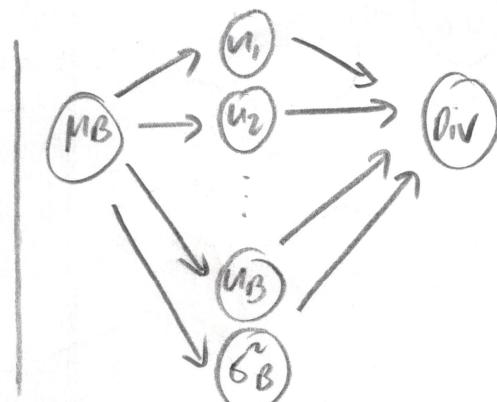
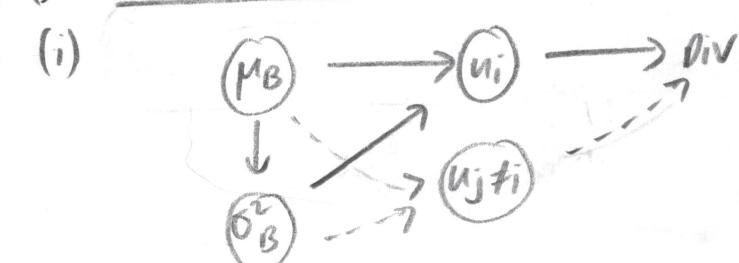
Hence  $\frac{\partial \text{div}}{\partial \sigma_B^2} = \sum_{i=1}^B \frac{\partial \text{div}}{\partial u_i} \frac{\partial u_i}{\partial \sigma_B^2} = \sum_{i=1}^B \frac{\partial \text{div}}{\partial u_i} \cdot -\frac{1}{2}(z_i - \mu_B)(\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$

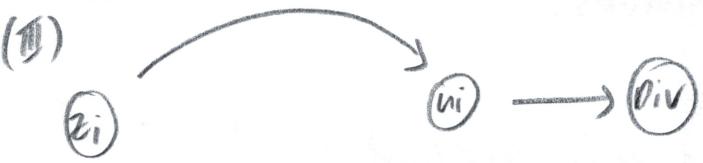
(\*)  $\rightarrow \frac{\partial \text{div}}{\partial \sigma_B^2} = -\frac{1}{2}(\sigma_B^2 + \epsilon)^{-\frac{3}{2}} \sum_{i=1}^B \frac{\partial \text{div}}{\partial u_i} (z_i - \mu_B)$

(II)(i)  $\frac{\partial \sigma_B^2}{\partial z_i} = \frac{2(z_i - \mu_B)}{B}$



### (II) Influence diagram (trajec.)



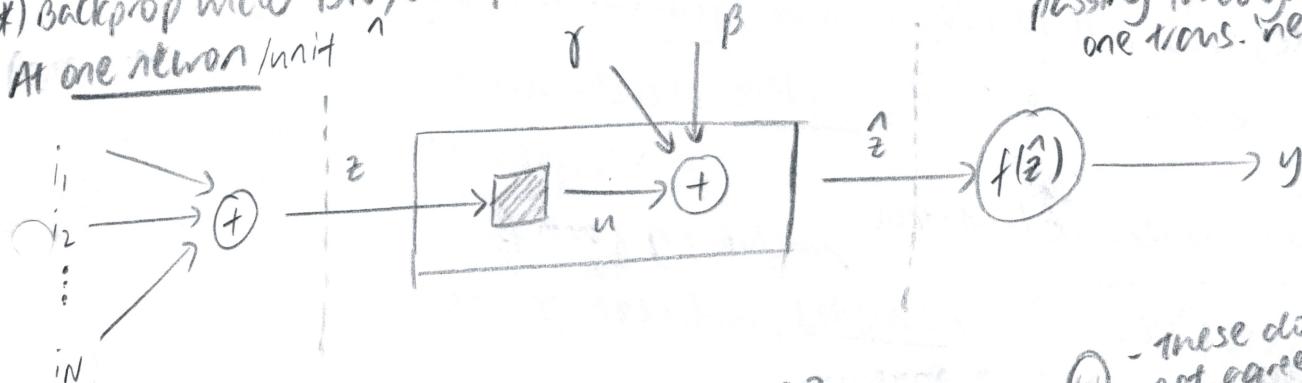


$$u_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

From above:-

$$\frac{\partial \text{Div}}{\partial u_i} \cdot \frac{\partial u_i}{\partial z_i} = \frac{\partial \text{Div}}{\partial u_i} \left( \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \right)$$

(\*) Backprop under BN, all eq. summarised:-  
At one neuron/unit



### FORWARD

$$z_i = \sum_j w_{ij} i_j + b$$

$$\mu_B = \frac{1}{B} \sum_{i=1}^B z_i$$

$$\sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (z_i - \mu_B)^2$$

$$u_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\hat{z}_i = y_i u_i + \beta$$

$$y = f(\hat{z}_i)$$

### BACKPROP

$$\frac{\partial \text{Div}}{\partial \hat{z}_i} = f'(\hat{z}_i) \frac{\partial \text{Div}}{\partial y_i}$$

$$\frac{\partial \text{Div}}{\partial y} = u_i \frac{\partial \text{Div}}{\partial \hat{z}_i}$$

$$\frac{\partial \text{Div}}{\partial \beta} = \frac{\partial \text{Div}}{\partial \hat{z}_i}$$

$$\frac{\partial \text{Div}}{\partial u_i} = \gamma \frac{\partial \text{Div}}{\partial \hat{z}_i}$$

$$\frac{\partial \text{Div}}{\partial z_i} = \frac{\partial \text{Div}}{\partial u_i} \frac{\partial u_i}{\partial z_i} + \frac{\partial \text{Div}}{\partial \sigma_B^2} \frac{\partial \sigma_B^2}{\partial z_i} + \frac{\partial \text{Div}}{\partial \mu_B} \frac{\partial \mu_B}{\partial z_i}$$

$$\frac{\partial \text{Div}}{\partial z_i} = \frac{\partial \text{Div}}{\partial u_i} \left( \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial \text{Div}}{\partial \sigma_B^2} \left( \frac{2(z_i - \mu_B)}{B} \right) + \frac{\partial \text{Div}}{\partial \mu_B} \cdot \frac{1}{B}$$

where  $\frac{\partial \text{Div}}{\partial \sigma_B^2} = -\frac{1}{2} (\sigma_B^2 + \epsilon)^{-3/2} \sum_{i=1}^B \frac{\partial \text{Div}}{\partial u_i} (z_i - \mu_B)$  ;  $\frac{\partial \text{Div}}{\partial \mu_B} = \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \sum_{i=1}^B \frac{\partial \text{Div}}{\partial u_i}$

(1)  
(2)

This only shows BN

transformations applied  
to one training instance  
passing through  
one trans. neuron.

(3)

These do  
not agree  
with Hui et al.,  
nor Ioffe &  
Szegedy (2015)

(\*) Backprop continues as usual from  $\frac{\partial \text{Div}}{\partial z_i}$  onwards.

⑩

⑪: for HW1; and completeness,  
it is possible to write the calculations in a more general form

- see supplementary notes with HW.

→ using a vector notation? (although time of writing; it is unclear  
why vector notation is used with summations rather than  
dot products; as in why not economise given more abstract rep)

- many cases; this will be required for efficient implementation in  
Numpy (\*) And above is not complete as a derivation... ("")

→ supplement → see HW1  
as well for  
implementation  
details.

(\*) Batch Norm - Inference / test stage

, including  $\gamma$  and  $\beta$

- after parameters have been estimated; we freeze 'parameters'  
- during test time; we require

$$\mu_{BN} = \frac{1}{N_B} \sum_{b=1}^{N_B} \mu_{B,b} \quad \text{where } N_B = \text{no. of batches}$$

$$\sigma_{BN}^2 = \frac{1}{(B-1)} \cdot \frac{1}{N_B} \sum_{b=1}^{N_B} \sigma_{B,b}^2$$

(\*) The key idea is that for an individual activation;  
we will have multiple computations of mini-batch means and  
variances  $\mu_B$  and  $\sigma_B^2$ .  
- more explicitly; if we have  $1, 2, 3, \dots, N_B$  mini-batches (i.e. partitions  
of the training set)

- we will have computed  $(\mu_{B,1}, \sigma_{B,1}^2)$  for batch 1

$$(\mu_{B,2}, \sigma_{B,2}^2) \text{ --- } 2$$

$$(\mu_{B,N_B}, \sigma_{B,N_B}^2) \text{ --- } N_B$$

- And that is for 1 activation unit

(\*) Some statistical reinforcement to link with Ioffe & Szegedy

- each mini-batch is a random sample of the entire training set (population); a mini-batch created by drawing  $B$  instances from training data.
- As data points comprising a mini-batch, randomly sampled, mini-batch statistics are fractions of randomly selected training instances (albeit deterministic)
- As mini-batch statistics can also be viewed as random variables; and as estimators of the mean and variance of each activation - Note  $z_1, z_2, \dots, z_N$

$$\rightarrow \mu_B = \frac{1}{B} \sum_{i=1}^B z_i \quad \sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (z_i - \mu_B)^2$$

- IID r.v.s.

(\*) To further clarify;  
we have a distri of the data  $f_z$

distri of statistic  $f_{\mu_B}$  and  $f_{\sigma_B^2}$

(\*)  $\mu_B$  and  $\sigma_B^2$  as mini-batch mean and variance are a little misleading from a stats perspective, <sup>notationally</sup>  $\mu_B$  should really be sample mean

$\bar{z}_B$ :  $\sigma_B^2$  is actually an estimator of the variance of that activation, with no bias adjustment (biased sample variance)

(\*) During inference; we use population statistics rather than mini-batch statistics.

- So we use the mini-batch sample means and sample variances to get our population statistics:

(or an estimate of them if implementation needs to that)

$$E[z_i] = \mu_{BN} = \frac{1}{NB} \sum_{b=1}^{NB} \mu_{B,b} = E_{\mu_B \sim p(\mu_B)} [\mu_B]$$

$$var[z_i] = \sigma_{BN}^2 = \frac{B}{B-1} \cdot \frac{1}{NB} \sum_{b=1}^{NB} \sigma_{B,b}^2 = \frac{B}{B-1} E_{\sigma_B^2 \sim p(\sigma_B^2)} [\sigma_B^2]$$

(unbiased variance estimate)

(\*) Note:-

$$\begin{aligned}\mathbb{E}[z_i] &= \frac{1}{NB} \sum_{b=1}^{NB} \mu_{B,b} = \frac{1}{NB} (\mu_{B,1} + \mu_{B,2} + \dots + \mu_{B,NB}) \\ &= \frac{1}{NB} \left( \frac{1}{B} \sum_{i=1}^B z_i + \frac{1}{B} \sum_{j=1}^B z_j + \dots + \frac{1}{B} \sum_{n=1}^B z_n \right) \\ &= \frac{1}{NBB} \sum_{i=1}^{NBB} z_i = \frac{1}{N} \sum_{i=1}^N z_i\end{aligned}$$

and similarly

$$\text{Var}[z_i] = \frac{B}{B-1} \mathbb{E}_{\delta_B^2} [\delta_B^2] = \mathbb{E} \left[ \frac{B}{B-1} \delta_B^2 \right] = \mathbb{E} [S_B^2]$$

→ relation between variance of activations under distri of data  
and the mean of the distribution of the unbiased sampling var of the activation.

(\*) implementation requires a running estimate being kept;  
and may prevent above from holding in some cases

① - once training complete:- for each activation and training example  
(after forward est.)



### 1) Normalisation

$$n_i = \frac{z_i - \mu_{BN}}{\sqrt{\sigma_{BN}^2 + \epsilon}}$$

(using population  
statistics NOT mini-batch  
stats)

### 2) linear transform :-

$$\hat{z}_i = \gamma \left( \frac{z_i - \mu_{BN}}{\sqrt{\sigma_{BN}^2 + \epsilon}} \right) + \beta = \frac{\gamma}{\sqrt{\sigma_{BN}^2 + \epsilon}} z_i + \left( \beta - \frac{\gamma \mu_{BN}}{\sqrt{\sigma_{BN}^2 + \epsilon}} \right)$$

(Ioffe & Szegedy  
2016)

(\*) see pseudocode of original paper,  
and also some calc / implement details on HW1.

$$\theta = \underset{\Theta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N L(x_i, \theta)$$

$\{x_1, \dots, x_N\}$  - training set

$\{x_1, \dots, x_m\}$  - minibatch

DL (for anchoring context) :-

$$J^*(\theta) = \sum_x \sum_y \text{Pdata}(x, y) L(f(x; \theta), y)$$

$$g = \nabla_{\theta} J^*(\theta) = \sum_x \sum_y \text{Pdata}(x, y) \nabla_{\theta} L(f(x; \theta), y)$$

$$\hat{g} = \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$$

use minibatch to approx. gradient of loss via:-

$$\frac{1}{m} \frac{\partial L(x_i; \theta)}{\partial \theta}$$

(\*) minibatch:

i) gradient our loss over minibatch is an estimate of loss gradient over entire training set.

ii) parallel computation.

(\*) rationale / observation for BN

-  $L = F_2(F_1(u, \theta_1), \theta_2)$        $F_1, F_2$  - arbitrary trans.

-  $L = F_2(x, \theta_2)$  where  $x = F_1(u, \theta_1)$

- gradient descent:-

$$\theta_2 \leftarrow \theta_2 - \frac{\alpha}{m} \frac{\partial F_2(x_i, \theta_2)}{\partial \theta_2}$$

- batch size  $m$   
- learning rate  $\alpha$

is equivalent to standalone network  $F_2$  with input  $x$ .

①  $\rightarrow$  not entirely clear.

$\Rightarrow$  keep distn of  $x$  fixed  $\Rightarrow \theta_2$  does not adjust to compensate for change in distn of  $x$

(\*) BN as altern. strategy for saturation, vanishing grad.  $\rightarrow$  spare time  
in saddle point

(f) Fixed distal inputs to a subnet; consequences outside subnet.  
 (sigmoid act.)  
 + weight matrix

$$z = g(Wu + b)$$

y-layer input      w-weight matrix  
 u-bias vector       $g(x) = \frac{1}{1 + \exp(-x)}$

- $|x| \rightarrow \infty$ ;  $g'(x) \rightarrow 0$
  - $\forall x_1, x_2, \dots, x_R$  in  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_K \end{pmatrix} = w_n + b$

- All aim except those with small abs. values will train slowly

- All dimensions except those with  $\Delta$  are affected by  $W$ ,  $D$  and parameters of layers below

- $\alpha$  unaffected by  $W$ ,  $b$  and parameters of layers ?
- $\Delta$  to those from during training  $\rightarrow$  "more many dimensions of  $\vec{x}$ "  
into the saturated regime of  
non-linearity, and slow converge

• Amplified as network depth ↑

- Amplified as network depth
- Known as saturation problem  $\Rightarrow$  vanishing gradients

(\*) Rectified with:-

i) ReLU (Nair & Hinton 2010)

- i) ReLU (Nair & Hinton 2010)
- ii) careful mit (Bergio & Glorot 2010; Saxe et al. 2013)

iii) small learning rates

- iii) small learning rates

(\*) Ensure that distn of nonlinearity inputs remains more stable as network trains then; optimiser will not get stuck in sat. regime and training would accelerate

## - key ideas

- key issues  
(\*) internal covariate shift - change in distri of internal nodes in network (network activations)

(\*) Batch Normalisation - reduce covariate shift; accelerate ANN training



## (\*) Key claims

- 1) Beneficial effect on gradient flow  
↳ reduce dependence of gradients on scale of param / initial vals.
- 2) Allows use of higher learning rates without risk of divergence
- 3) regularisation
- 4) reduces need for dropout (Srivastava et al. 2014)

## 2. Discussion of trajectory → reducing ICS

### 3. Norm. via mini-batch stats

(\*) Normalise each scalar feature  $\text{muf}$ .  
- i.e. zero mean, unit variance

TRAIN SET

(\*) Layer with  $d$ -dim input:-

$$\underline{x} = \begin{pmatrix} x^{(1)} \\ \vdots \\ x^{(d)} \end{pmatrix}$$

(\*) Normalise each dim :-

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}} \quad \forall k=1, \dots, d \quad \text{over training set}$$

(\*) Introduce for each activation  $x^{(k)}$ , a pair of param.  $\gamma^{(k)}$  and  $\beta^{(k)}$   
which scale and shift normalised value  $\hat{x}^{(k)}$ :-

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)} \quad \forall k=1, \dots, d$$

(\*) Estimate  $\gamma^{(k)}$  and  $\beta^{(k)}$  ( $\forall k=1, \dots, d$ ) with orig param.

(\*) Reduce rep. para of net; compare  $x^{(k)}$  (original act.) with

$$\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]} \quad \text{and} \quad \beta^{(k)} = E[x^{(k)}]$$

(\*) Accy & for mini-batches

(\*) minibatch size  $m$ .

(\*) focus on single activ./dim  $x^{(k)}$ , drop  $k$ .

(\*)  $m$ -values of activation in mini-batch

$$B = \{x_1, x_2, \dots, x_m\} = \{\hat{x}_1 \dots m\}$$

(\*) normalised values  $\{\hat{x}_1 \dots m\}$

(\*) linear trans  $\{y_1 \dots m\}$

(\*) BN trans:  $BN_{\gamma, \beta}: x_1 \dots m \rightarrow y_1 \dots m$

- BN transform Algo (one activ./dim) NOT instance

IN: values of  $x$  over mini-batch  $B = \{x_1 \dots m\}$

OUT:  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

c-num stab const

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta = BN_{\gamma, \beta}(x_i)$$

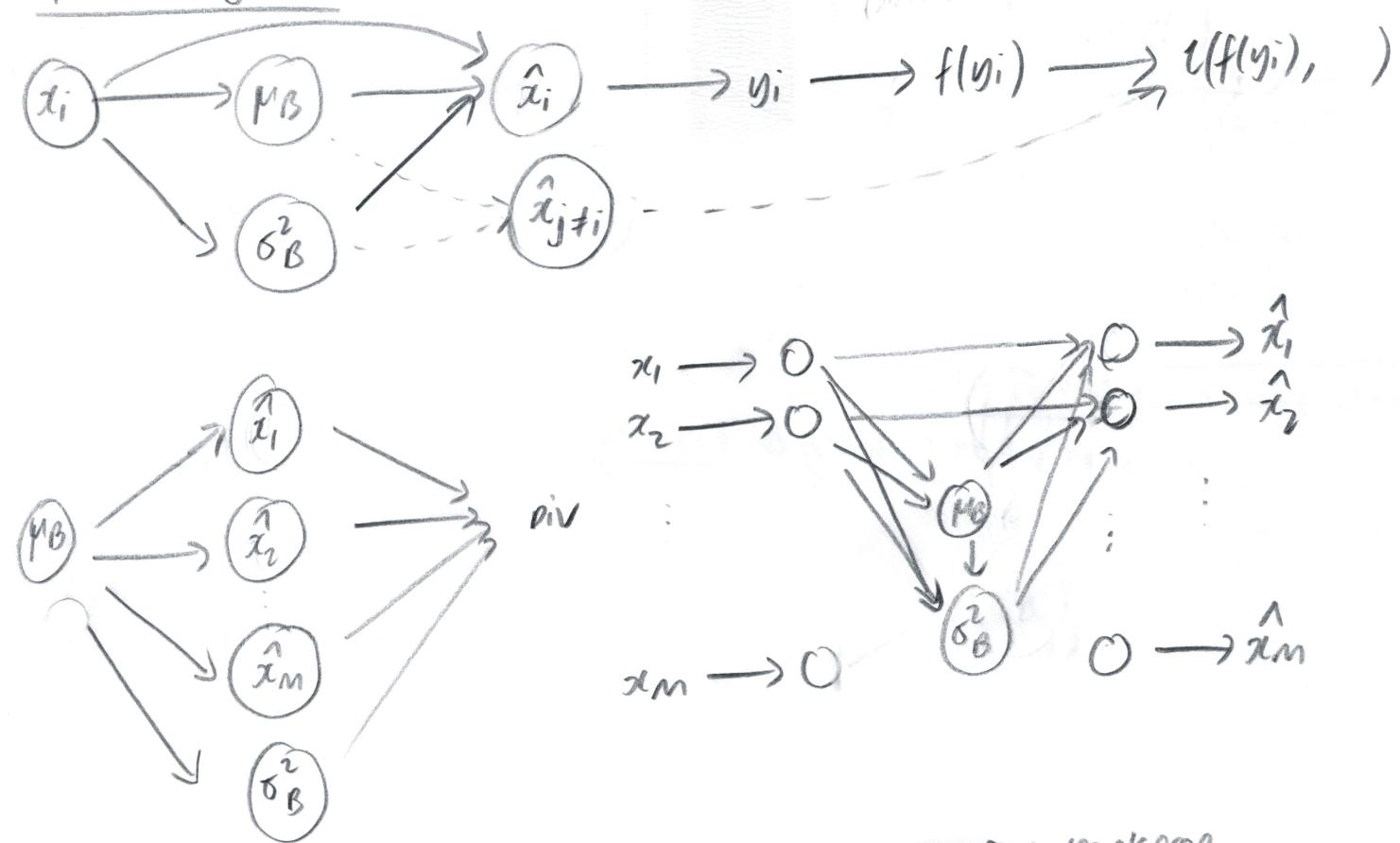
(\*) Just because we consider one mini/batch  $\Rightarrow$  scalar fn.

(\*)  $BN_{\gamma, \beta}(x_i) = f(x_i, x_{-i}, \mu_B(x_1 \dots m), \sigma_B^2(x_1 \dots m), \gamma, \beta)$   
incorrect (?)

(\*) lectures  $\rightarrow$  batch norm vector  
from our minibatch.

# Backpropagation (BN)

get adept  
at not being conf. by art.



(W): I don't fully understand the form of the chain rule for backprop in BN

(W): I can't fully visualise the influence chain; in particular why the lectures' divergence omits the late activations from the influence diagram.

• (W) How does this affect backprop algo with respect to  $\nabla_{w_R} CE(\hat{y}, t)$  and  $\nabla_{b_R} CE(\hat{y}, t)$

(\*) write down anyway (Ioffe & Szegedy) (need to reconcile) with rec.

$$\frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \cdot \frac{\partial y_i}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \cdot \gamma$$

$$\frac{\partial l}{\partial \sigma_B^2} = \left( \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \cdot \frac{\sum_{i=1}^m -2(x_i - \mu_B)}{m}$$

$$\frac{\partial \underline{L}}{\partial \mu_B} = \frac{\partial \underline{L}}{\partial x_i} = \frac{\partial \underline{L}}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial \underline{L}}{\partial \sigma_B} \cdot \frac{2(x_i - \mu_B)}{m} + \frac{\partial \underline{L}}{\partial \mu_B} \cdot \frac{1}{m}$$

$$\frac{\partial \underline{L}}{\partial y} = \sum_{i=1}^m \frac{\partial \underline{L}}{\partial y_i} \cdot \hat{x}_i$$

$$\frac{\partial \underline{L}}{\partial \beta} = \sum_{i=1}^m \frac{\partial \underline{L}}{\partial y_i}$$

(\*) These have been simplified

(\*) Need to work out what's going on exactly