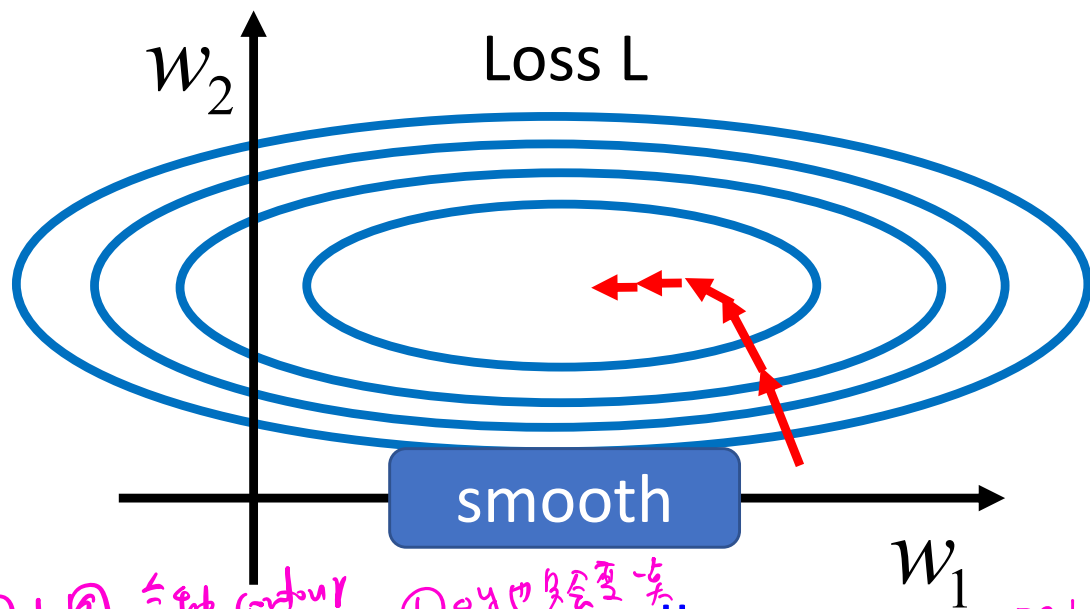


這是講 training 的 tips
... 等 2 的 CNN, seq2seq 都會用到
... 現在才講

Quick Introduction of Batch Normalization

Hung-yi Lee 李宏毅

Changing Landscape



① error surface 很崎岖的时候
 比较难 train, 那我们能不能
 直接把山搓平, 让他变的
 比较好 train?

⇒ batch normalization 就是一种方法

② 这边即使不是 convex 的问题
 也不见得好 train,

如图, 当 w_1 较缓 (走很大步, loss 才降一点)
 w_2 较陡 (走一小步, loss 就降很多)

用固定的 learning rate
 也不行,
 所以有各种
 调 lr 的方式
 e.g. Adam, momentum

④ 如果 x_1 的值都很小

那现在用力一点解它

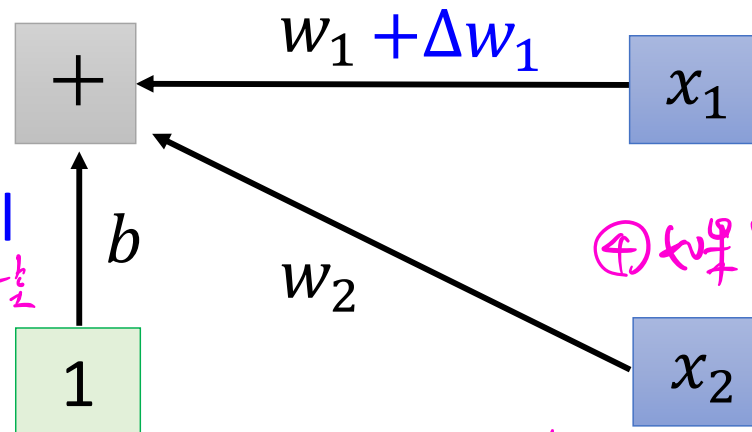
③ 上图这种 contour
 和 x_1, x_2 的 scale
 有关

⑥ y 也会变 $small$
 $+ \Delta y$

⑤ 那 w_1 变慢一点

$$L = \sum e$$

$+ \Delta L$
 small

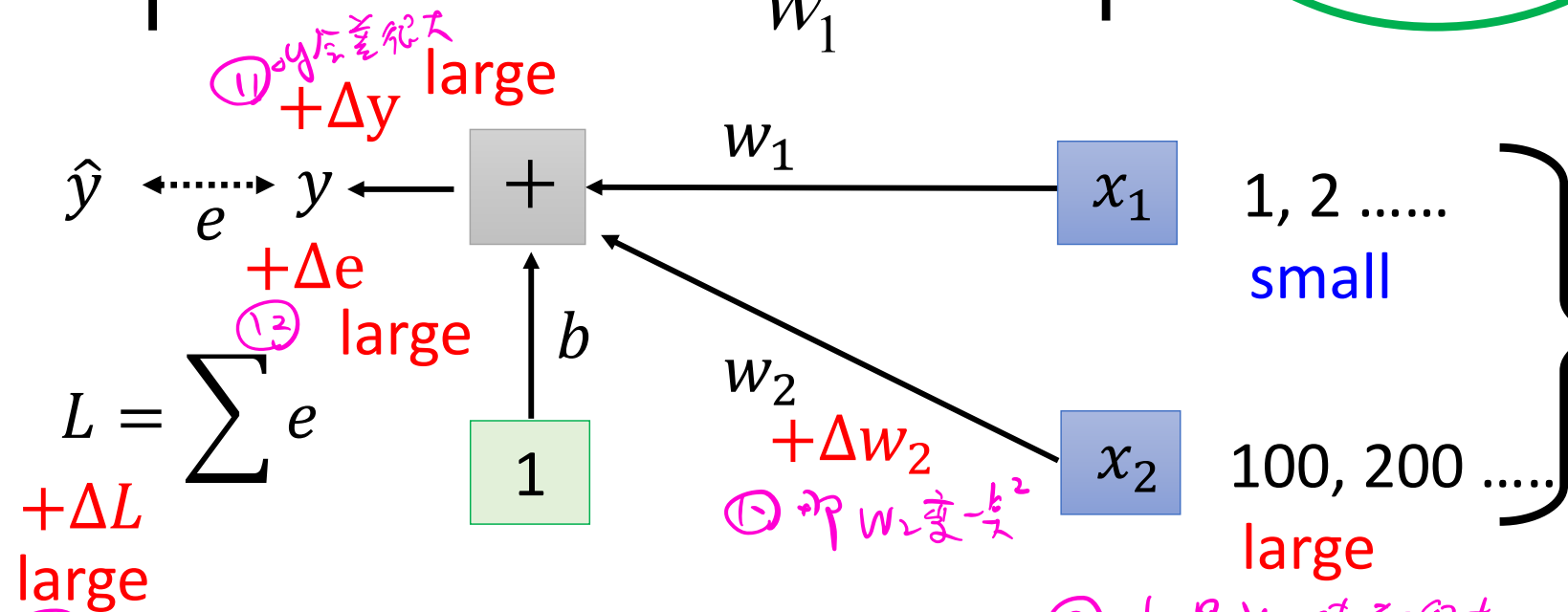
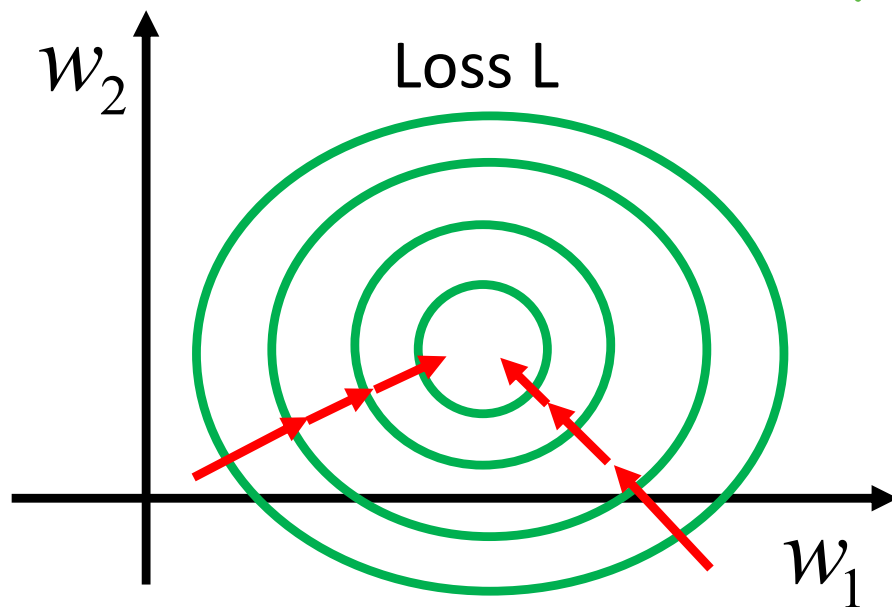
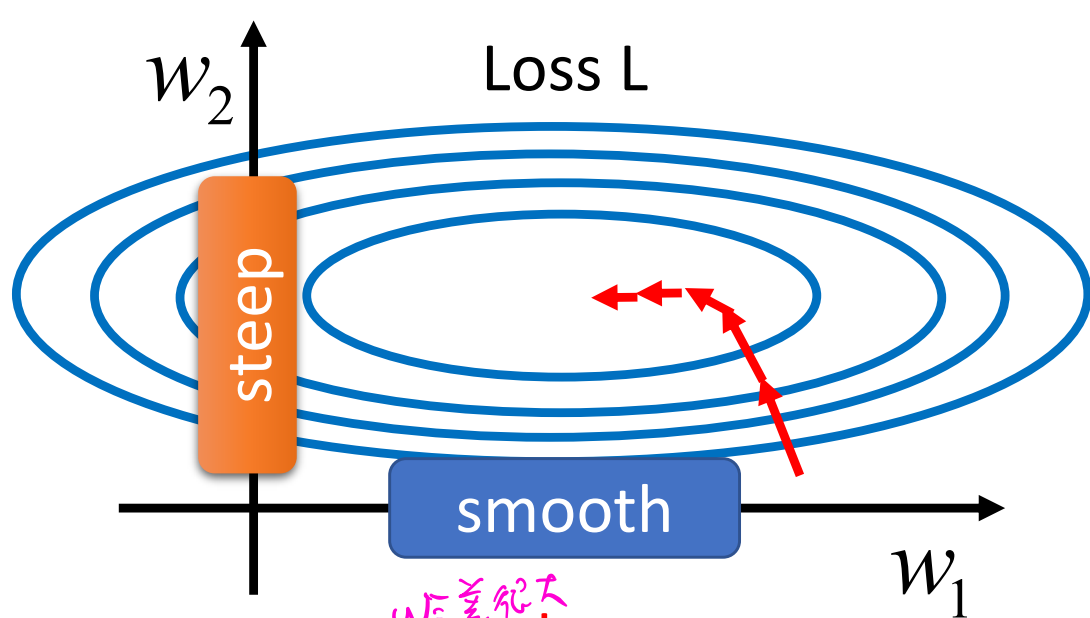


1, 2
 small

⑦ $loss$ 是 $loss^2$ 就是上图很缓的样子

Changing Landscape

⑭ 直接把 error surface 改成 那即使是非凸的 gradient descent 也能好 train 起来



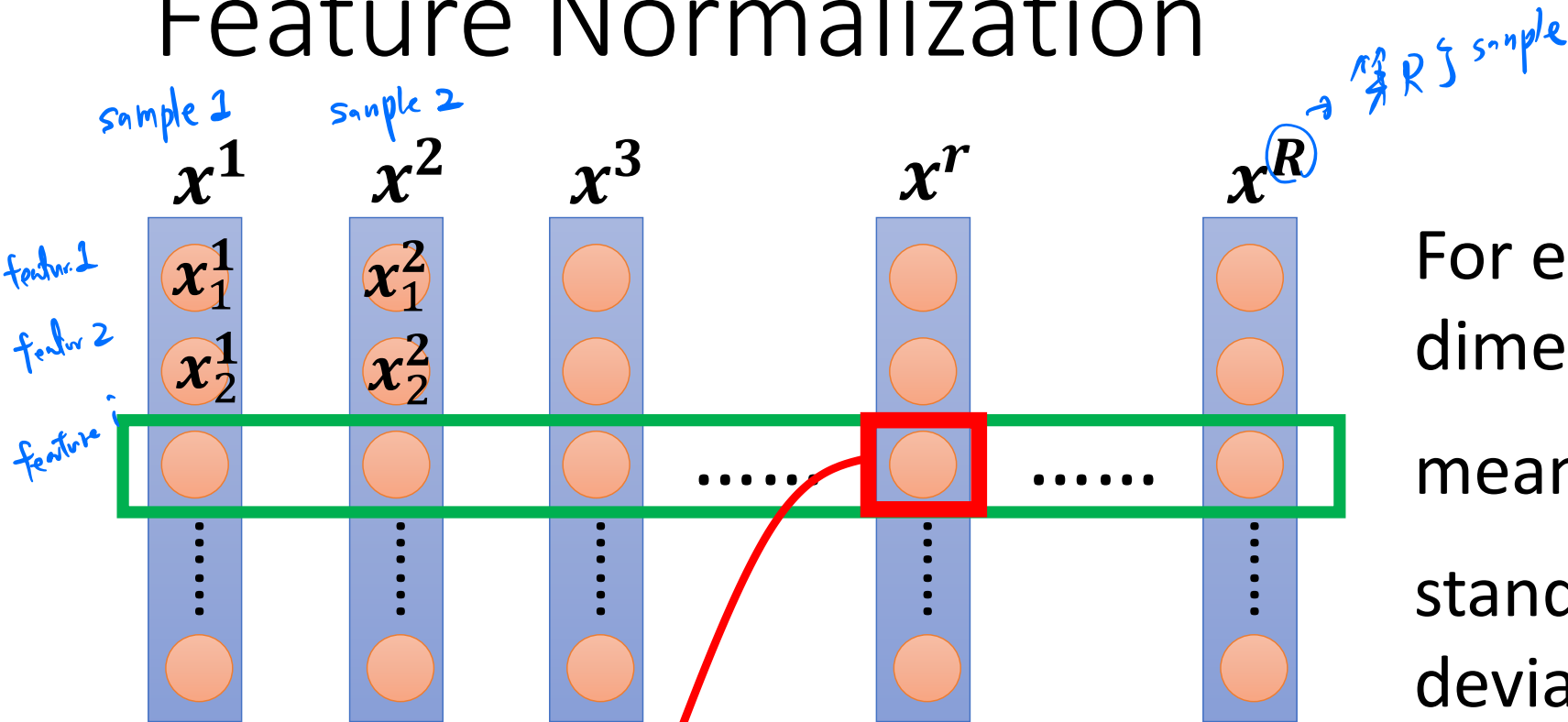
same range

⑮ 做层 很简单, 3 层

⑰ 如果 x_2 值即很大

x_1, x_2
变-样本的
scale
更好

Feature Normalization



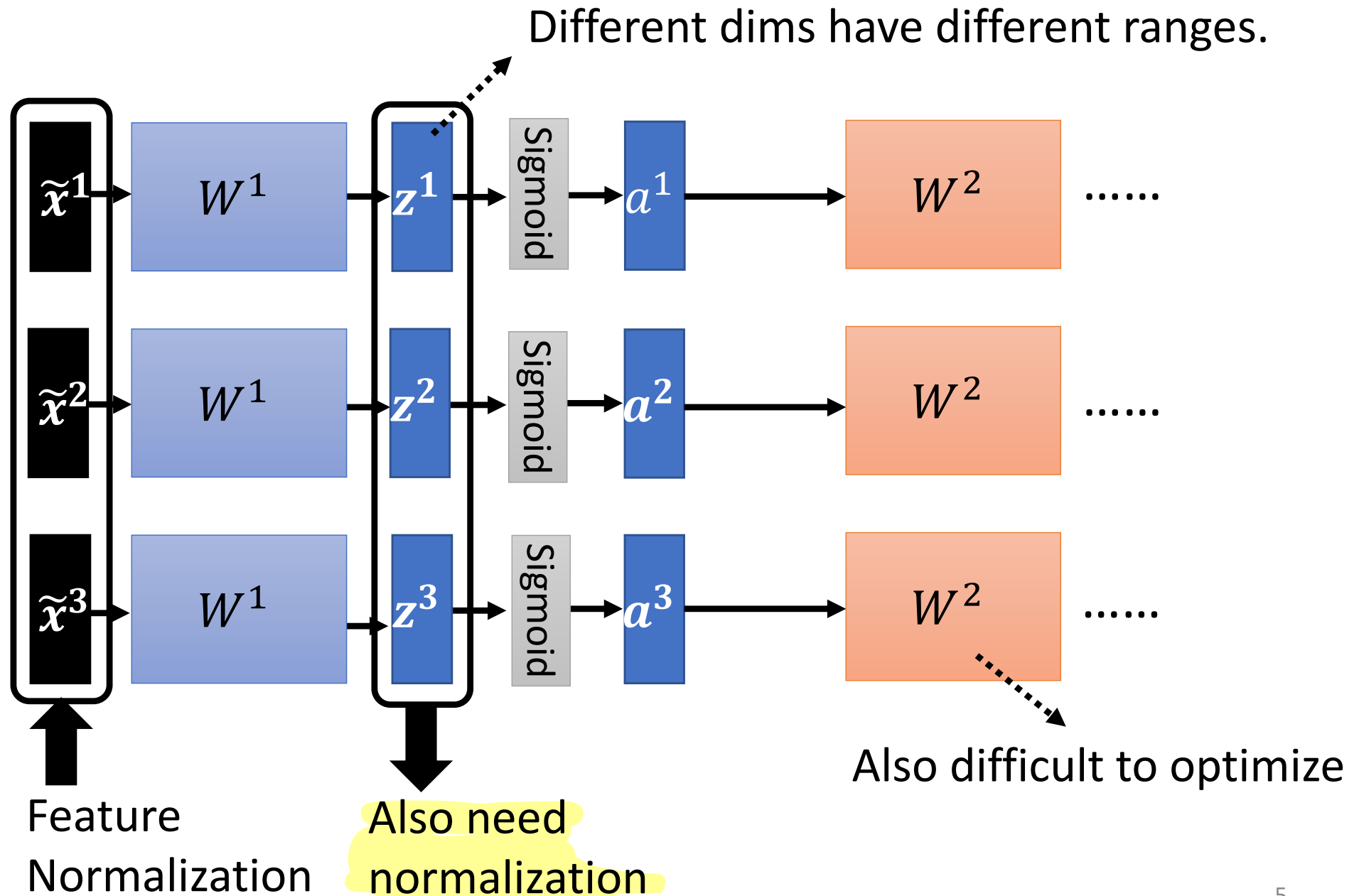
For each dimension i :
mean: m_i
standard deviation: σ_i

$$\tilde{x}_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

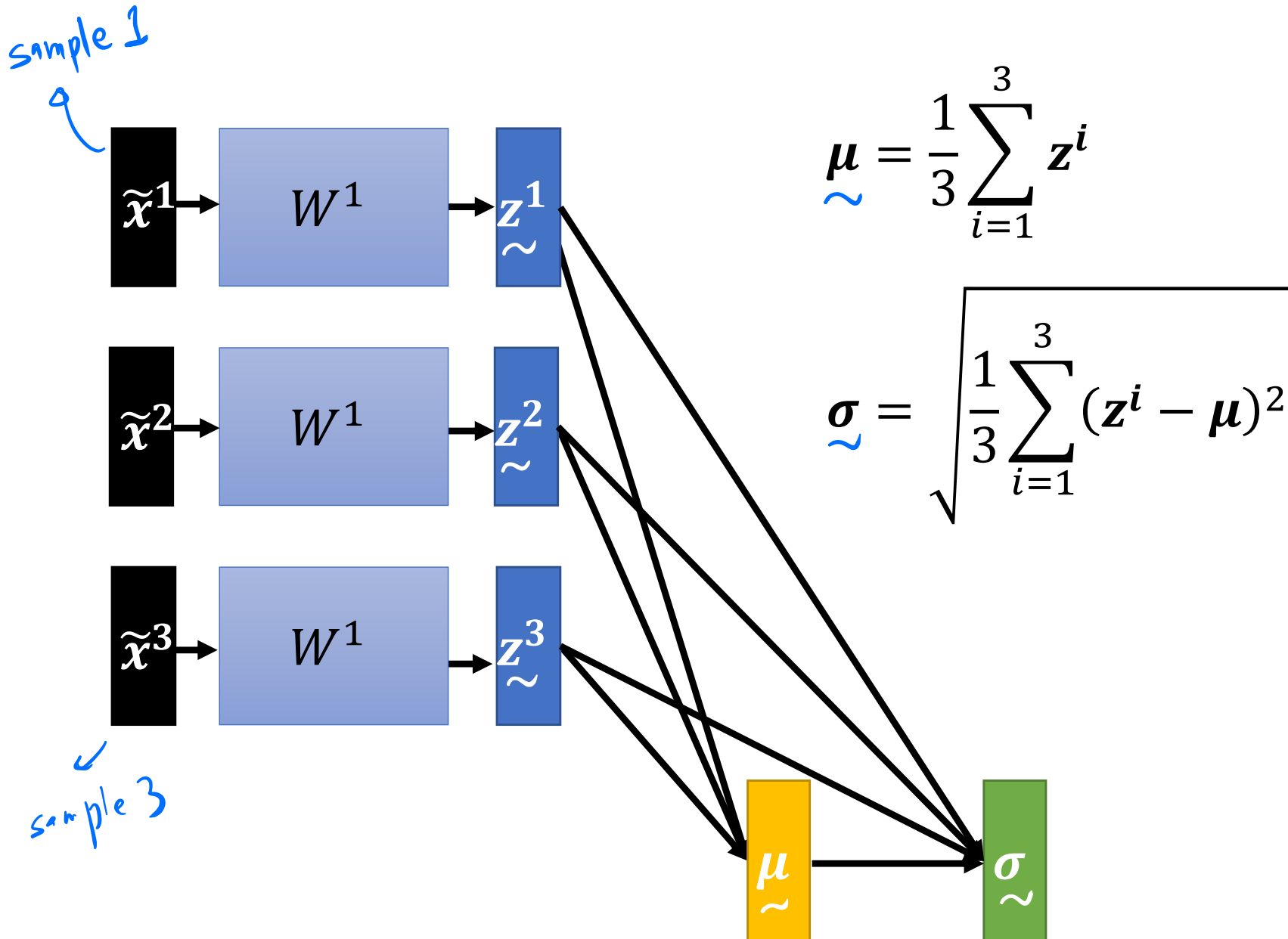
The means of all dims are 0,
and the variances are all 1

In general, feature normalization makes gradient descent converge faster.

Considering Deep Learning



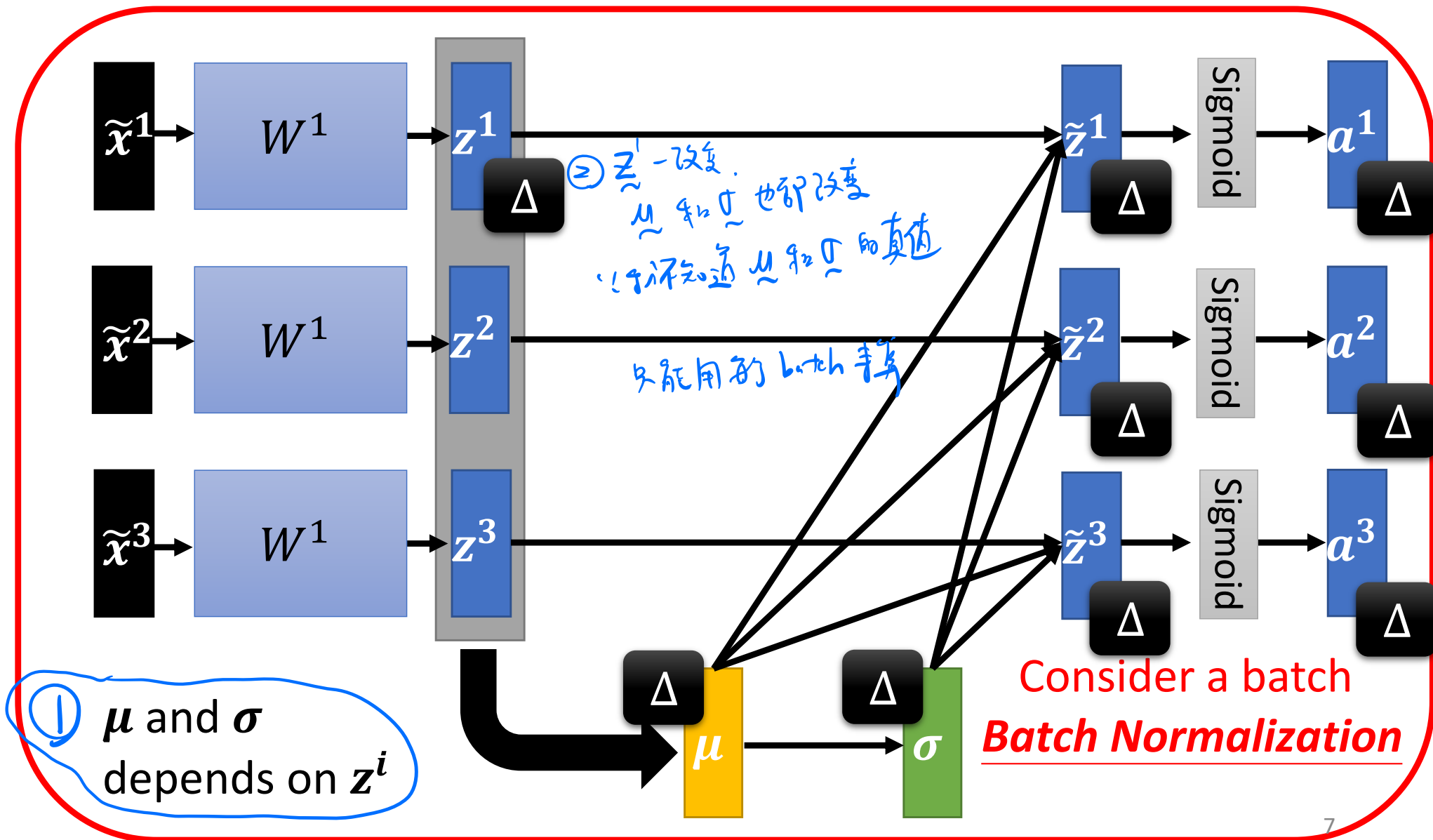
Considering Deep Learning



Considering Deep Learning

$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

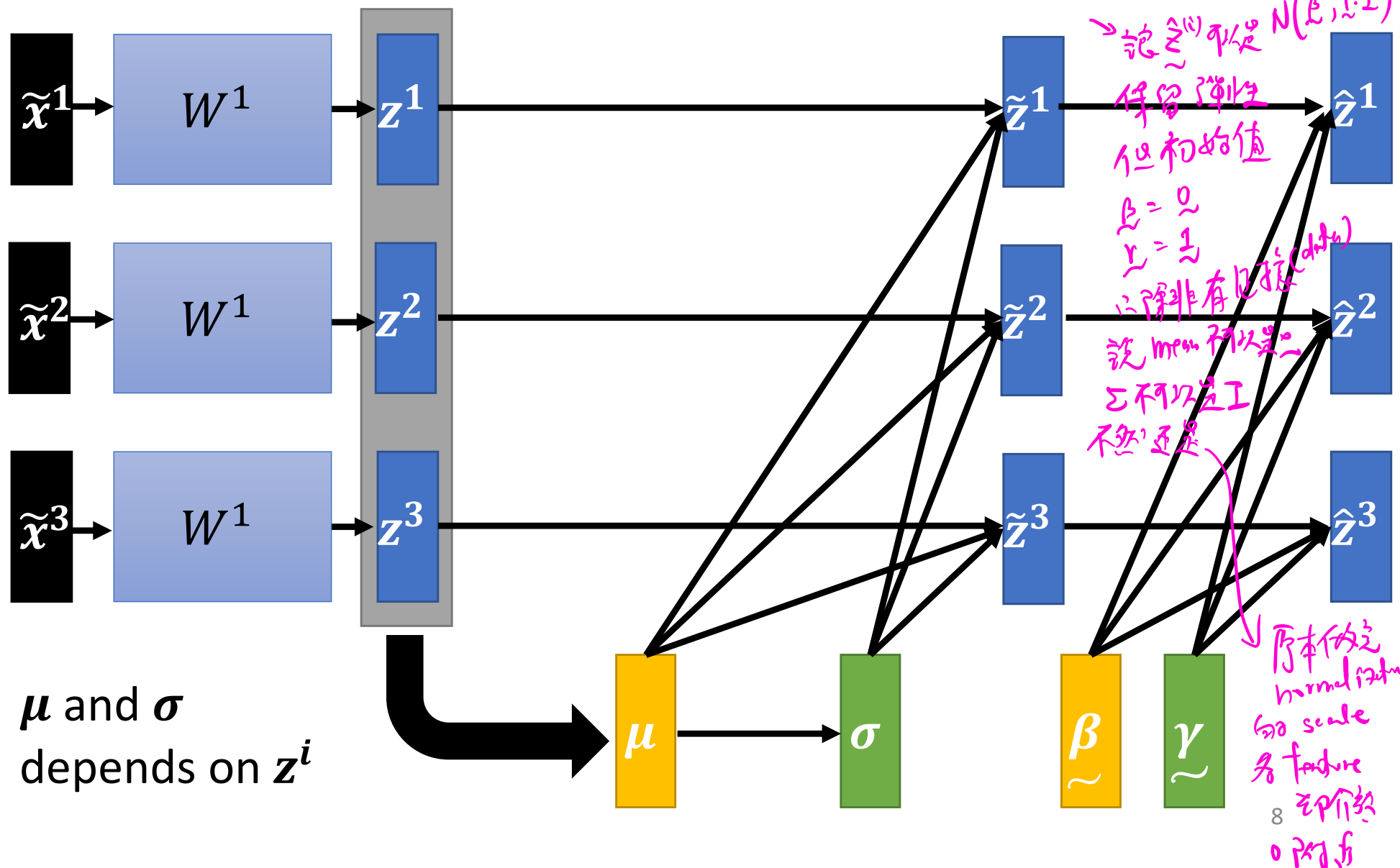
This is a large network!



Batch normalization

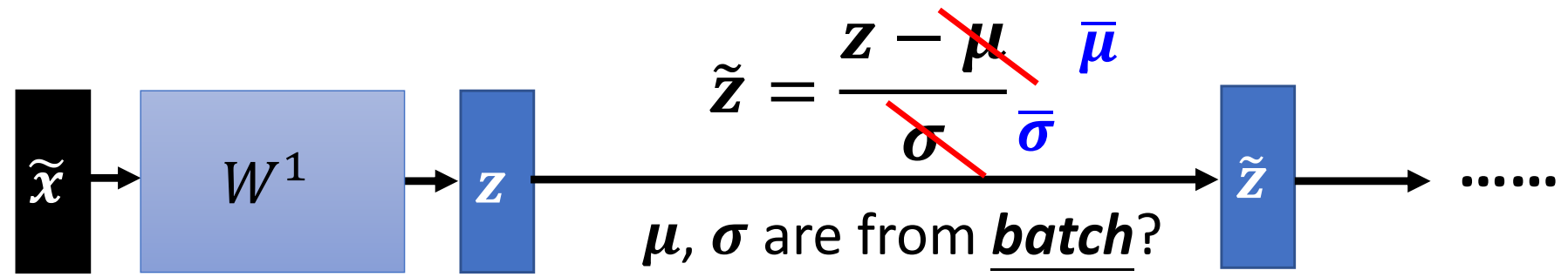
$$\tilde{z}^i = \frac{z^i - \mu}{\sigma} \quad \tilde{z}^i \sim N(0, 1)$$

$$\hat{z}^i = \gamma \odot \tilde{z}^i + \beta$$



Batch normalization – Testing

上阶段, μ 和 σ 是用 training 时得到的 $\bar{\mu}$ 和 $\bar{\sigma}$



We do not always have batch at testing stage.

Computing the moving average of μ and σ of the batches during training.

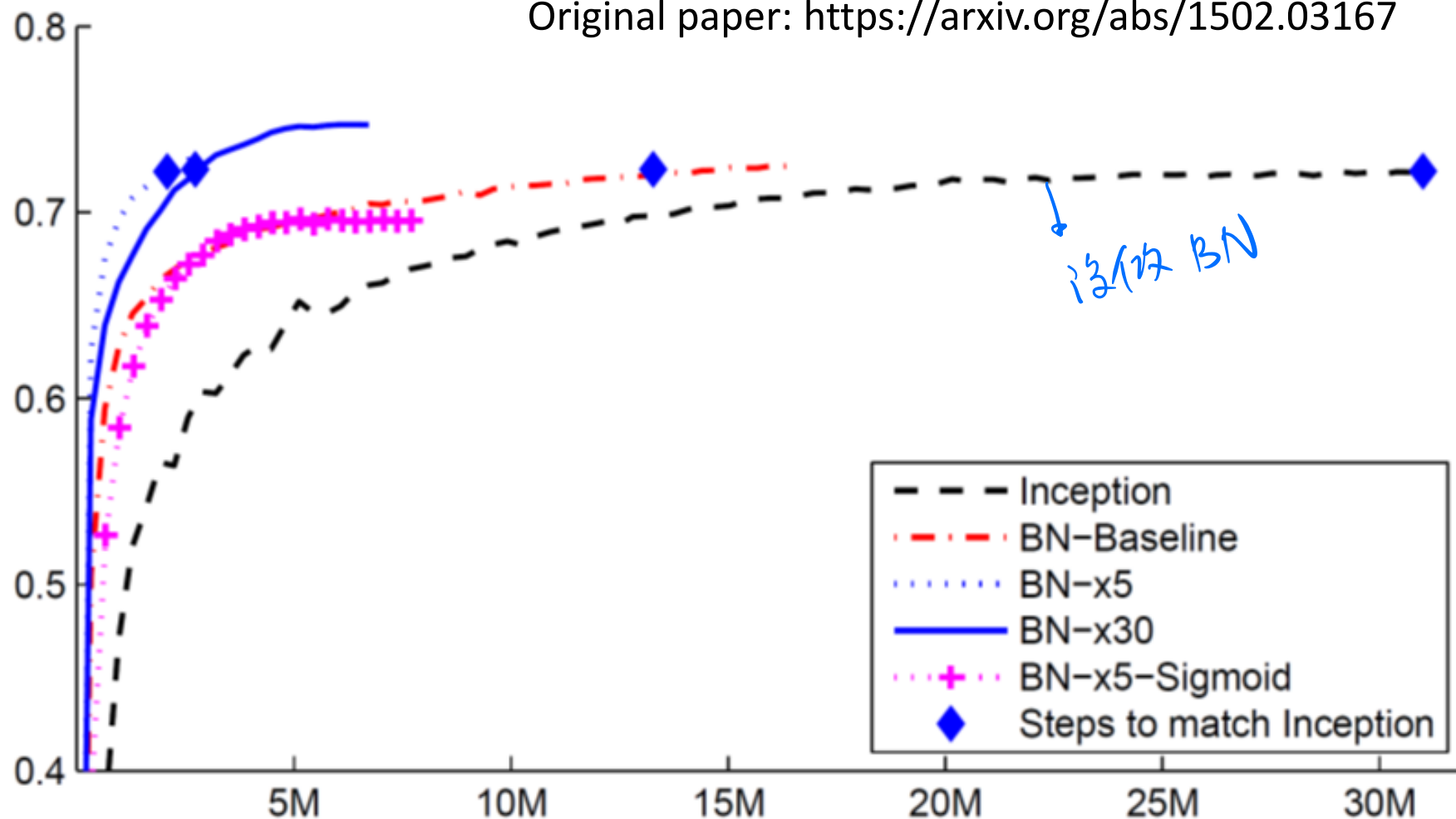
$$\mu^1 \quad \mu^2 \quad \mu^3 \quad \dots \quad \mu^t$$

$$\bar{\mu} \leftarrow p\bar{\mu} + (1 - p)\mu^t$$

Diagram showing the update of the moving average $\bar{\mu}$. The term $p\bar{\mu}$ is highlighted in yellow. A blue arrow points from the text "moving-average" to the $p\bar{\mu}$ term. Another blue arrow points from the text "batch" to the μ^t term.

Batch normalization

Original paper: <https://arxiv.org/abs/1502.03167>



hyperparameter
学习率 0.1
... 主要是在学习率上
做文章

没做 BN

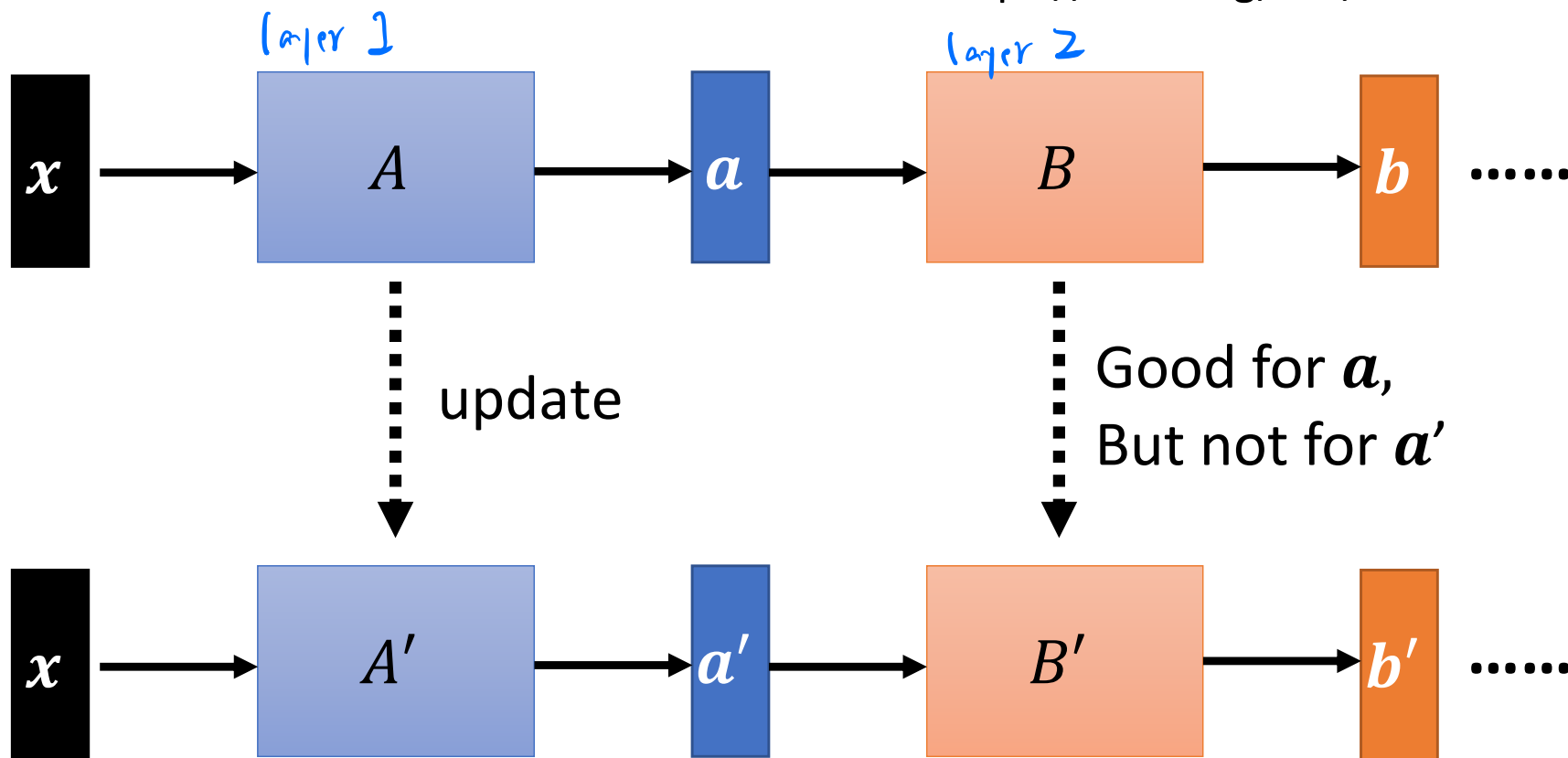
① 因为 BN 啥都有用，
在 gradient descent 中，
分布会收敛到更窄的分布

有研究嘗試證明BN能幫助優化，而提出的猜想，他說，因為在hidden layer (i.e. a)，BN能幫助優化，因為back propagation 沒有受到影響，所以效果不好，如果有的話，batch normalization 那 a 和 a' 的分布會一樣，就能training faster

Internal Covariate Shift?

How Does Batch Normalization Help Optimization?

<https://arxiv.org/abs/1805.11604>



Batch normalization make a and a' have similar statistics.

② Experimental results do not support the above idea.

Internal Covariate Shift?

How Does Batch Normalization Help Optimization?

<https://arxiv.org/abs/1805.11604>

③ BN 的代法, 大家还是认为是在 所得到的数据

Experimental results (and theoretically analysis) support batch normalization change the landscape of error surface.

and 12 of Appendix B.) This suggests that the positive impact of BatchNorm on training might be somewhat serendipitous. Therefore, it might be valuable to perform a principled exploration of the design space of normalization schemes as it can lead to better performance.

serendipitous (偶然的)

意料之外的发现

penicillin



To learn more

- Batch Renormalization
 - <https://arxiv.org/abs/1702.03275>
- Layer Normalization
 - <https://arxiv.org/abs/1607.06450>
- Instance Normalization
 - <https://arxiv.org/abs/1607.08022>
- Group Normalization
 - <https://arxiv.org/abs/1803.08494>
- Weight Normalization
 - <https://arxiv.org/abs/1602.07868>
- Spectrum Normalization
 - <https://arxiv.org/abs/1705.10941>

还有一项 normalization 的方法
可参考

(反正就是 training tips
有人试过这样做
normalization, 有效
他就提出来)

