

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

Hyperparameter tuning

Tuning process

Hyperparameters

→ α

β 0.9

$\beta_1, \beta_2, \epsilon$
0.9 0.999 10^{-8}

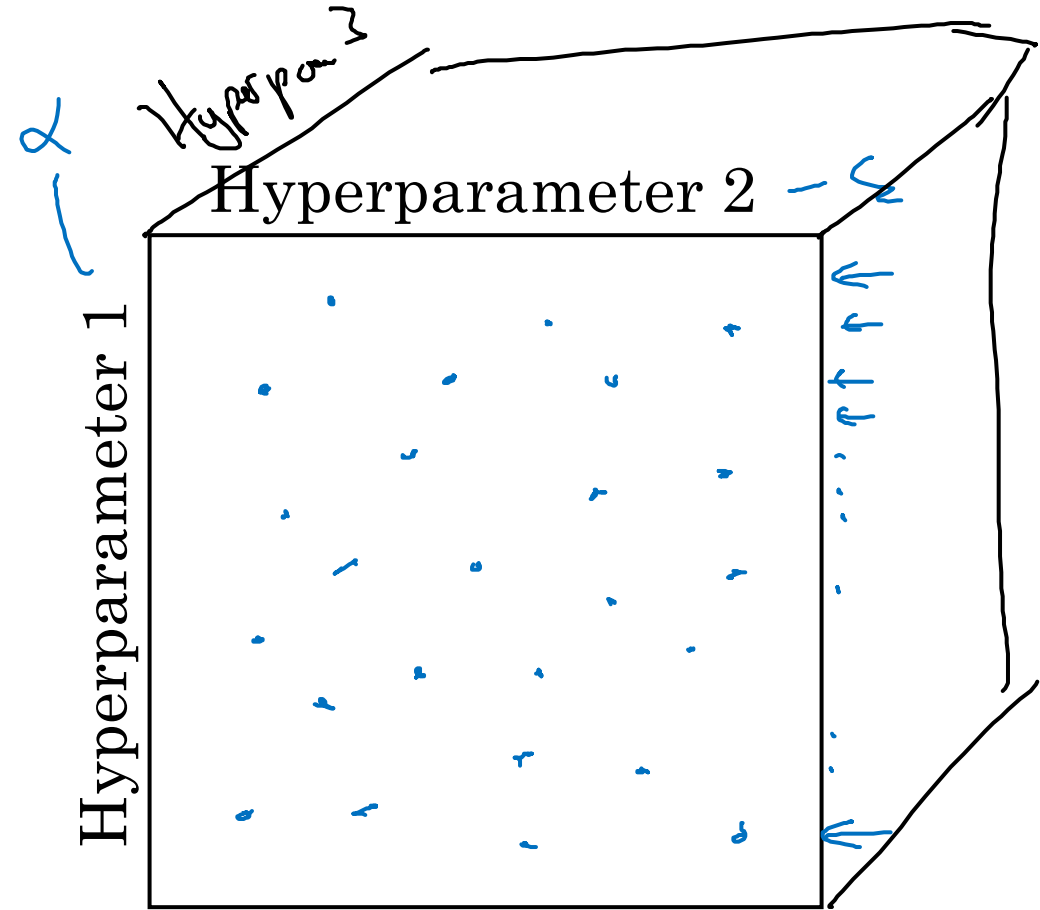
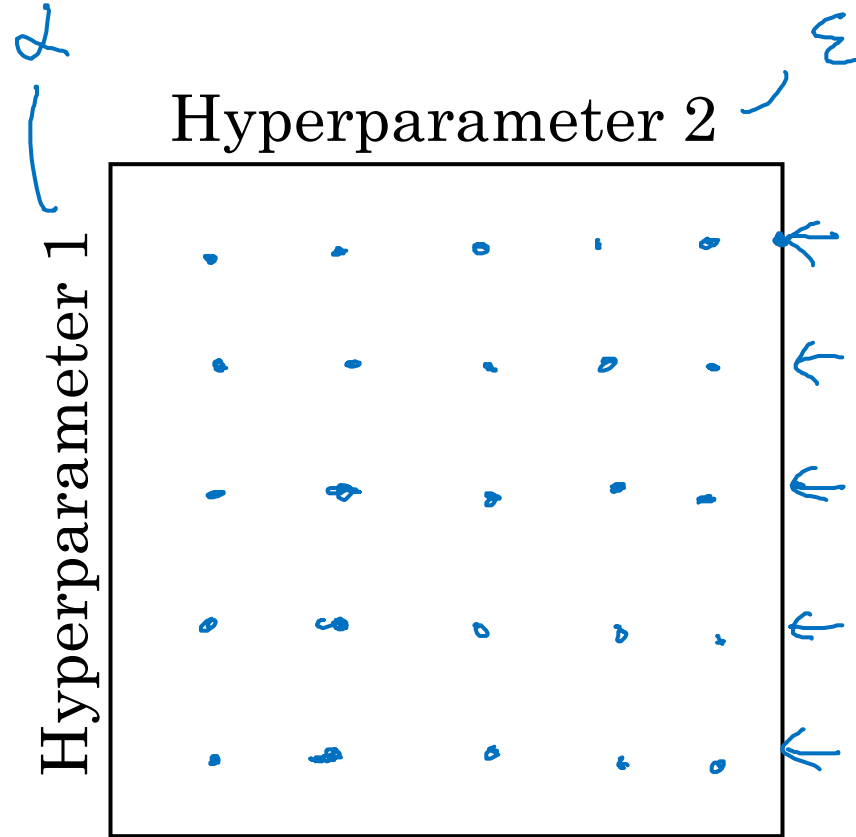
layers

hidden units

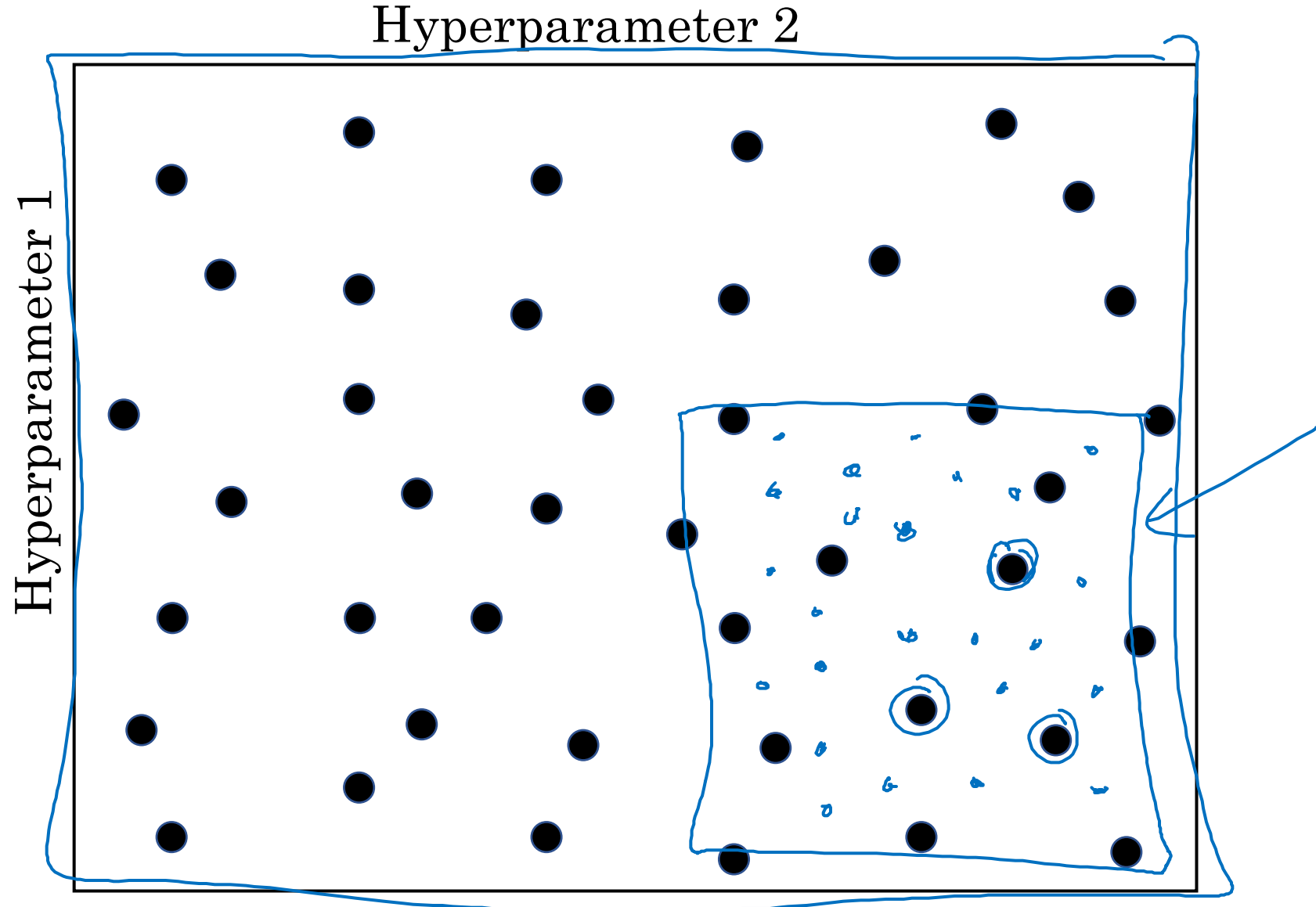
learning rate decay

mini-batch size

Try random values: Don't use a grid



Coarse to fine





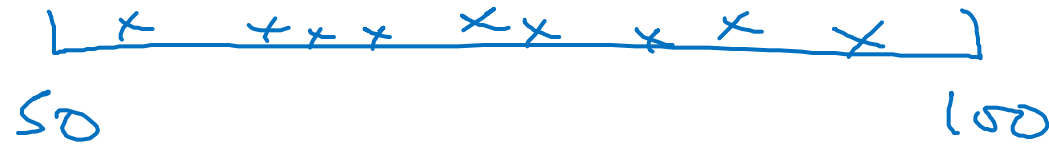
deeplearning.ai

Hyperparameter tuning

Using an appropriate
scale to pick
hyperparameters

Picking hyperparameters at random

→ $n^{\text{test}} = 50, \dots, 100$

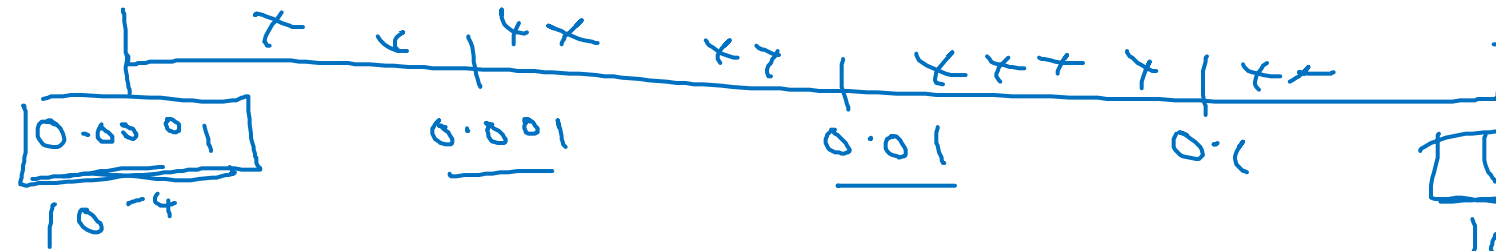
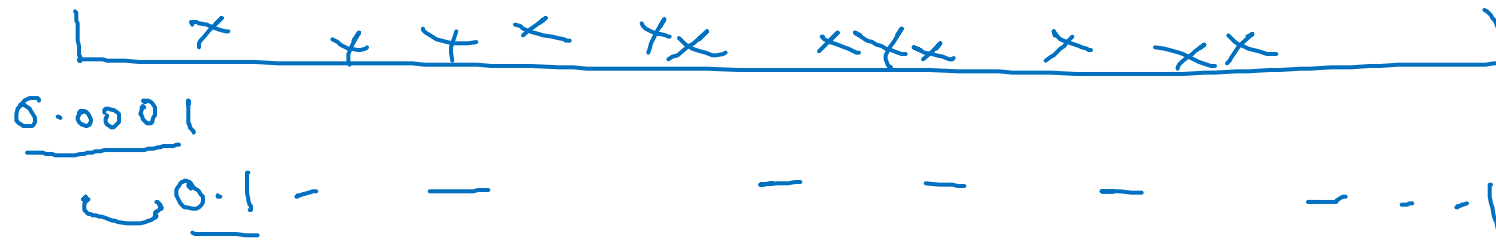


→ #layers $L : 2 - 4$

2, 3, 4

Appropriate scale for hyperparameters

$$\alpha = 0.0001, \dots, 1$$



$$10^a$$

$$a = \log_{10} 0.0001$$

$$= -4$$

$$r = -4 * \text{np.random.rand}()$$

$$\alpha = 10^r$$

$$\frac{10^a \dots 10^b}{}$$

$$\frac{r \in [a, b]}{[-4, 0]}$$

$$\underline{\alpha = 10^r}$$

$$\leftarrow r \in [-4, 0]$$

$$\leftarrow 10^{-4} \dots 10^0$$

$$\frac{b = \log_{10} 1}{= 0}$$

Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \quad \dots \quad 0.999$$

\downarrow
 10

\downarrow
 1000

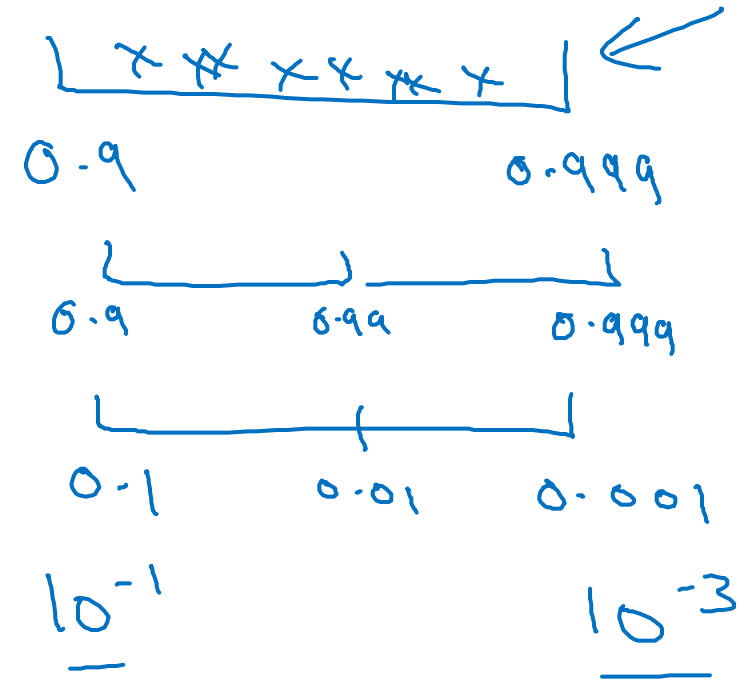
$$1 - \beta = 0.1 \quad \dots \quad 0.001$$

$$\beta: 0.999 \rightarrow 0.9995 \quad \} \sim 10$$

$$\beta: 0.999 \rightarrow 0.9995$$

~ 1000
 ~ 2000

$$\frac{1}{1 - \beta_K}$$



$$r \in [-3, -1]$$

$$1 - \beta = 10^r$$

$$\beta = 1 - 10^r$$

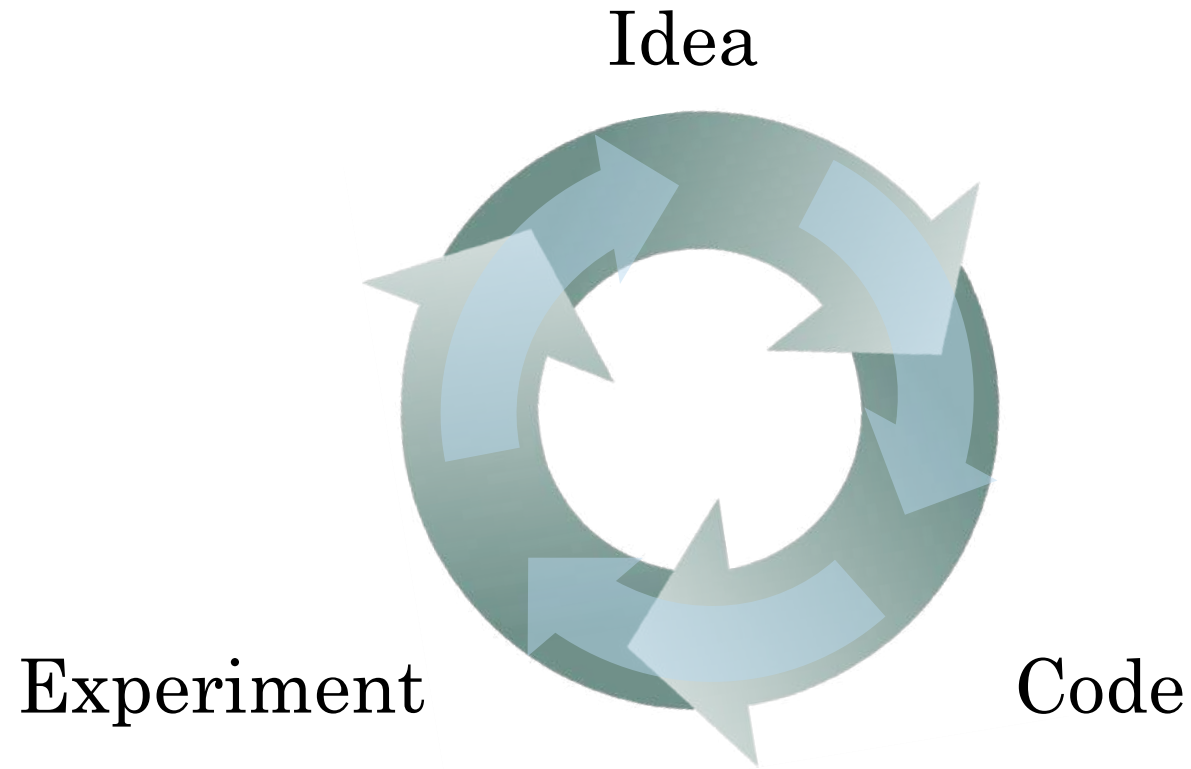


deeplearning.ai

Hyperparameters tuning

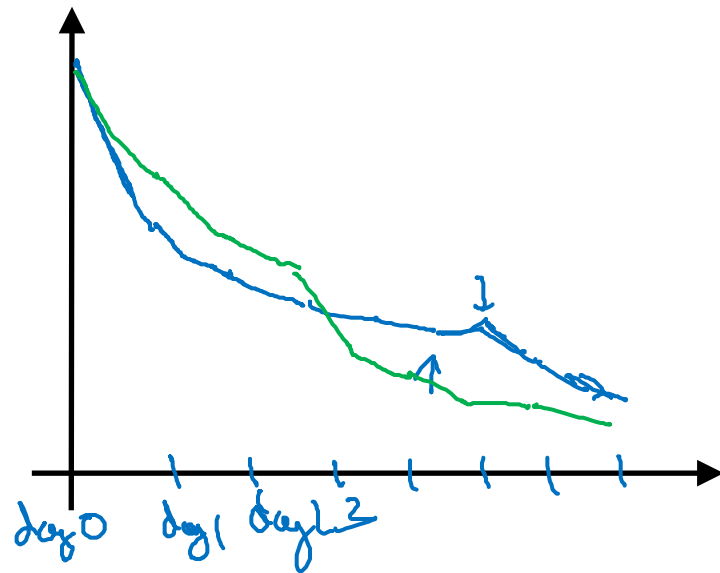
Hyperparameters
tuning in practice:
Pandas vs. Caviar

Re-test hyperparameters occasionally



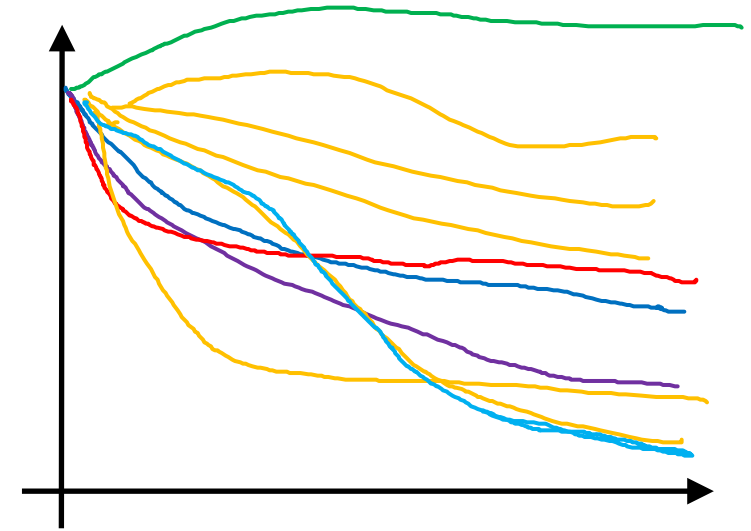
- NLP, Vision, Speech,
Ads, logistics,
- Intuitions do get stale.
Re-evaluate occasionally.

Babysitting one model



Panda ↵

Training many models in parallel



Caviar ↵

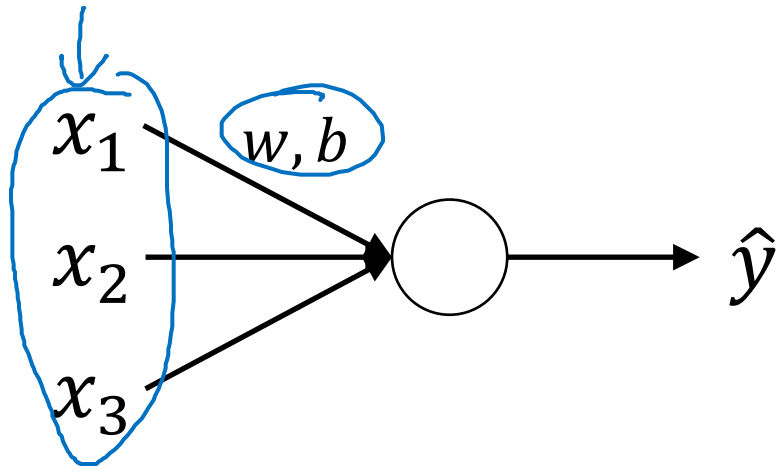


deeplearning.ai

Batch Normalization

Normalizing activations
in a network

Normalizing inputs to speed up learning

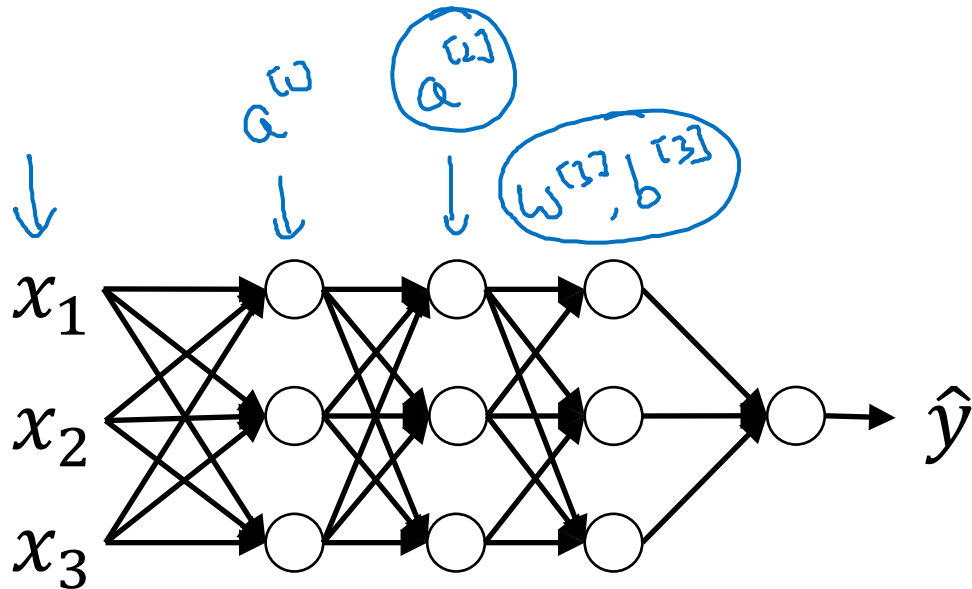
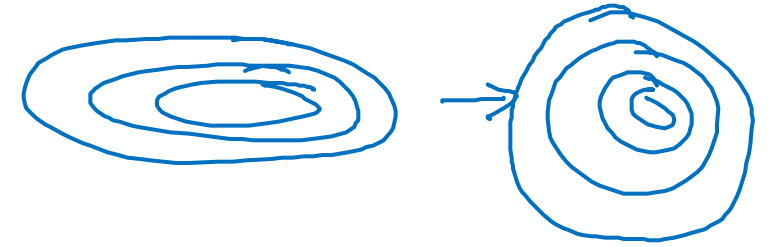


$$\mu = \frac{1}{n} \sum_i x^{(i)}$$

$$X = X - \mu$$

$$\sigma^2 = \frac{1}{n} \sum_i x^{(i)2} \quad \leftarrow \text{element-wise}$$

$$X = X / \sigma^2$$



Can we normalize $\frac{a^{[2]}}{w^{[2]}, b^{[2]}}$ so as to train faster

Normalize $\frac{z^{[2]}}{\uparrow}$

Implementing Batch Norm

Given some intermediate values in NN

$z^{(1)}, \dots, z^{(m)}$
 $z^{[l]}(i)$

$$\begin{aligned} \mu &= \frac{1}{m} \sum_i z^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \\ z_{\text{norm}}^{(i)} &= \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ \hat{z}^{(i)} &= \gamma z_{\text{norm}}^{(i)} + \beta \end{aligned}$$

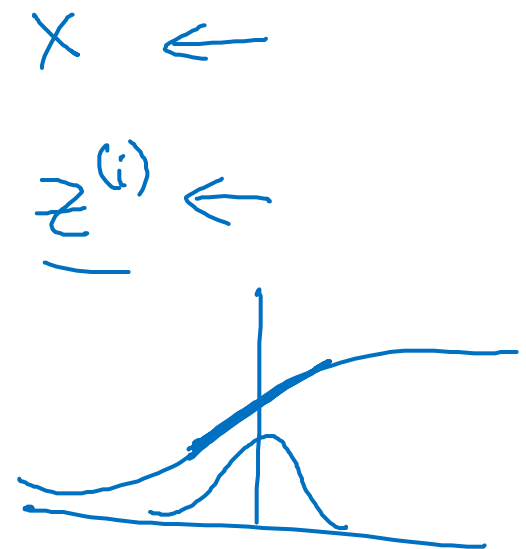
If

$$\gamma = \sqrt{\sigma^2 + \epsilon}$$

$$\beta = \mu$$

then $\hat{z}^{(i)} = z^{(i)}$

learnable parameters of model.



Use $\hat{z}^{[l]}(i)$ instead of $z^{[l]}(i)$.

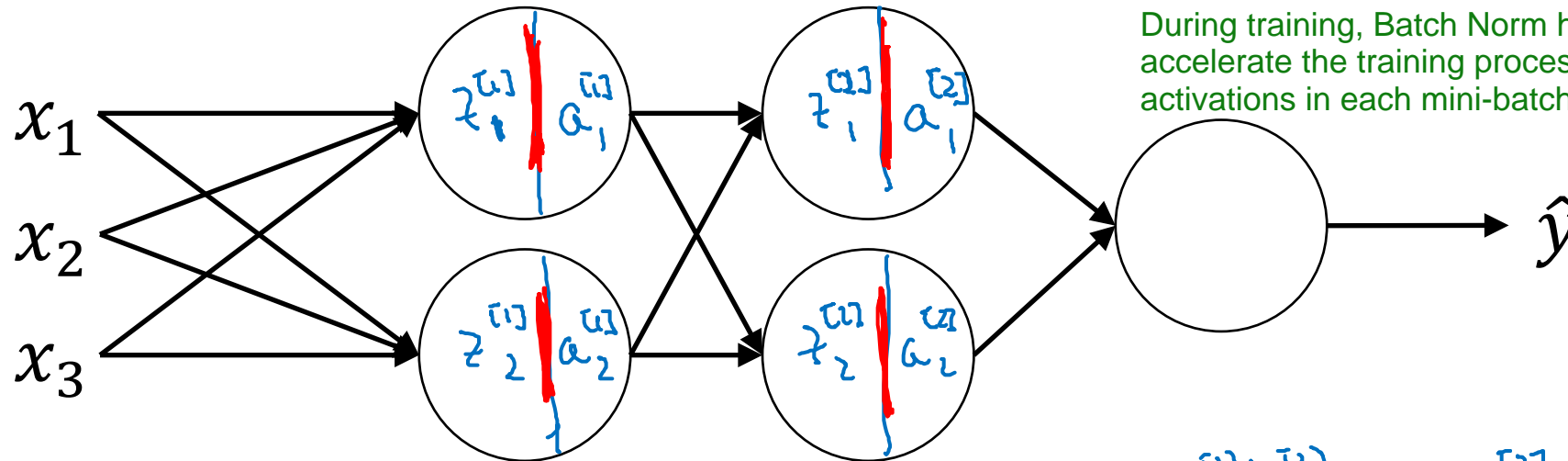


deeplearning.ai

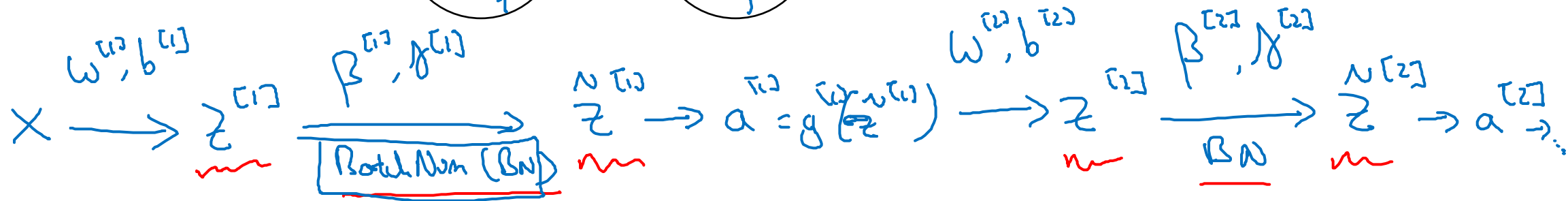
Batch Normalization

Fitting Batch Norm
into a neural network

Adding Batch Norm to a network



During training, Batch Norm helps stabilize and accelerate the training process by normalizing the activations in each mini-batch.



Parameters: $\left\{ W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots, W^{(L)}, b^{(L)} \right\}$
 $\rightarrow \underline{\beta}^{(1)}, \underline{\gamma}^{(1)}, \underline{\beta}^{(2)}, \underline{\gamma}^{(2)}, \dots, \underline{\beta}^{(L)}, \underline{\gamma}^{(L)}$
 $\rightarrow \underline{\beta}$

Gamma and Beta are added parameters in this list

$$d\beta^{(2)} \quad \beta = \beta - \alpha d\beta^{(2)}$$

tf.nn.batch-normalization ←

To implement batch norm in tensorflow

Working with mini-batches

$$\underline{X^{[1]}} \xrightarrow{W^{[1]}, b^{[1]}} \underline{Z^{[1]}} \xrightarrow[\text{BN}]{\beta^{[1]}, \gamma^{[1]}} \underline{\tilde{Z}^{[1]}} \rightarrow g^{[1]}(\tilde{Z}^{[1]}) = a^{[1]} \xrightarrow{W^{[2]}, b^{[2]}} \underline{Z^{[2]}} \rightarrow \dots$$

$$\boxed{X^{[2]}} \rightarrow \underline{Z^{[2]}} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \underline{\tilde{Z}^{[2]}} \rightarrow \dots$$

$$X^{[L]} \rightarrow \dots$$

During normalization, we will be subtracting mean of $Z[l]$ from $Z[l]$, but since all values of $b[l]$ will remain same hence it will eventually cancel out, so no need to store the values of $b[l]$

Parameters: $W^{[l]}, \cancel{b^{[l]}}, \beta^{[l]}, \gamma^{[l]}$

\downarrow \downarrow \downarrow
 $(n^{[l-1]}, 1)$ $(n^{[l-1]}, 1)$ $(n^{[l-1]}, 1)$

\uparrow

$Z^{[l-1]}_{(n^{[l-1]}, 1)}$

$$\rightarrow \underline{Z^{[l]}} = W^{[l]} a^{[l-1]} + \cancel{\boxed{b^{[l]}}}$$

$$Z^{[l]} = W^{[l]} a^{[l-1]}$$

Or you may think of it as :
Since in the normalization process the values of $Z[l]$ are re-centered at the origin, it is irrelevant to add the $b[l]$ parameter.

$$\rightarrow \underline{\tilde{Z}^{[l]}} = \gamma^{[l]} Z^{[l]}_{\text{norm}} + \boxed{\beta^{[l]}}$$

Implementing gradient descent

for $t = 1 \dots \text{num Mini Batches}$

Compute forward pass on $X^{\{t\}}$.

In each hidden layer, use BN to ~~compute~~ ^{REPLACE} $\underline{z}^{\{t\}}$ with $\underline{\hat{z}}^{\{t\}}$.

Use backprop to compute $\underline{dw}^{\{t\}}$, ~~$\underline{db}^{\{t\}}$~~ , $\underline{dp}^{\{t\}}$, $\underline{d\delta}^{\{t\}}$

Update params
$$\left. \begin{aligned} w^{\{t\}} &:= w^{\{t-1\}} - \alpha dw^{\{t\}} \\ \beta^{\{t\}} &:= \beta^{\{t-1\}} - \alpha dp^{\{t\}} \\ \gamma^{\{t\}} &:= \dots \end{aligned} \right\} \leftarrow$$

Works w/ momentum, RMSprop, Adam.

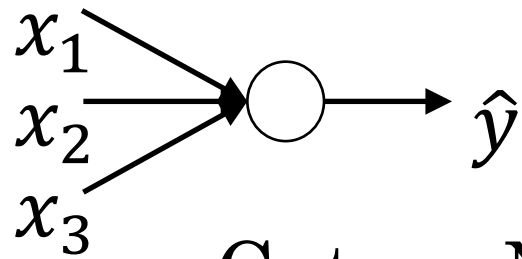


deeplearning.ai

Batch Normalization

Why does
Batch Norm work?

Learning on shifting input distribution

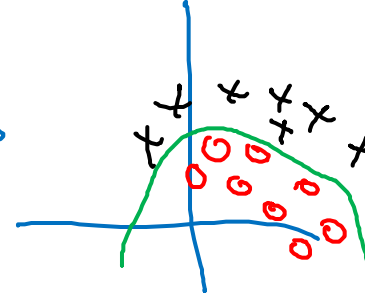
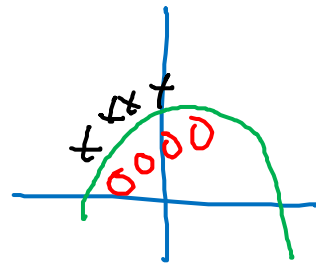


Cat

Non-Cat

$y = 1$

$y = 0$



It reduces internal covariate shift, making the network more robust and allowing for faster convergence.

The mean and variance statistics are computed per mini-batch, which introduces some randomness.



"Covariate shift"

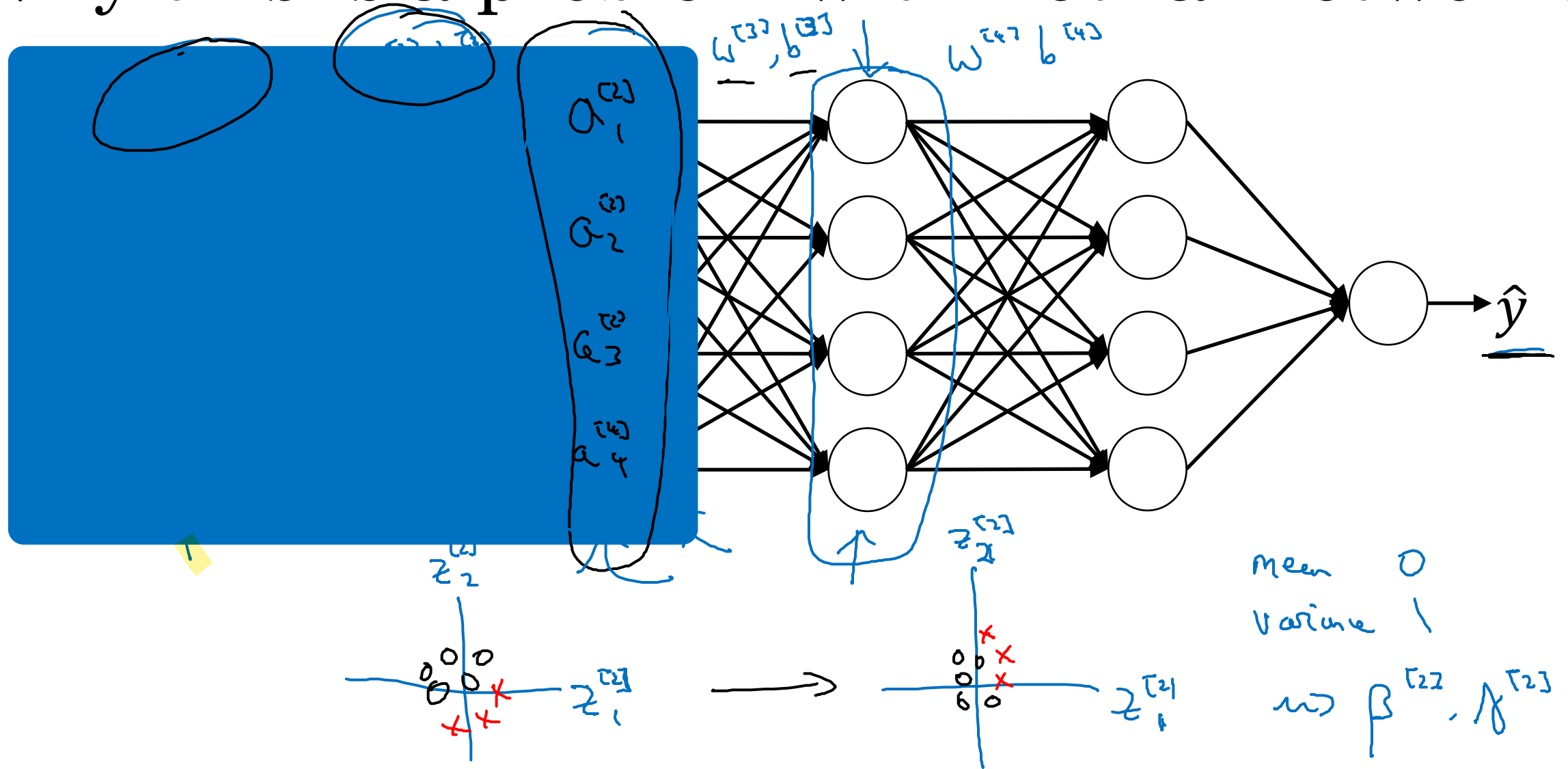
$\underline{x} \rightarrow y$

$y = 1$

$y = 0$



Why this is a problem with neural networks?



Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

Test Phase (Inference):

- When evaluating the neural network on unseen data (test/validation set or during deployment), we want consistent and deterministic predictions.
- If Batch Norm were applied during inference, it would use the mini-batch statistics from the test data, leading to different normalization parameters than those used during training.

Solution :

During the test, the parameters μ and σ^2 are estimated using an exponentially weighted average across mini-batches used during training.



deeplearning.ai

Multi-class classification

Softmax regression

Recognizing cats, dogs, and baby chicks



3



1



2



0



3



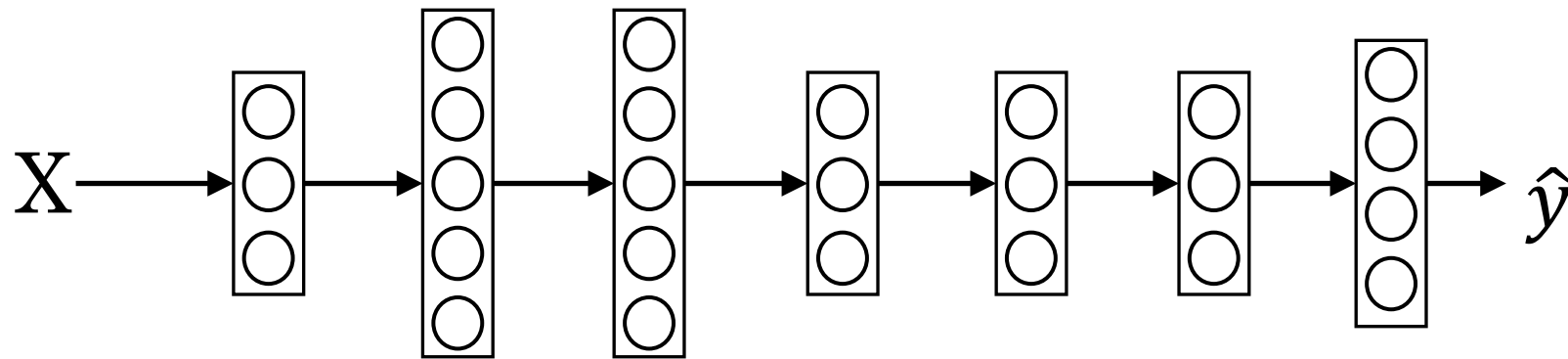
2



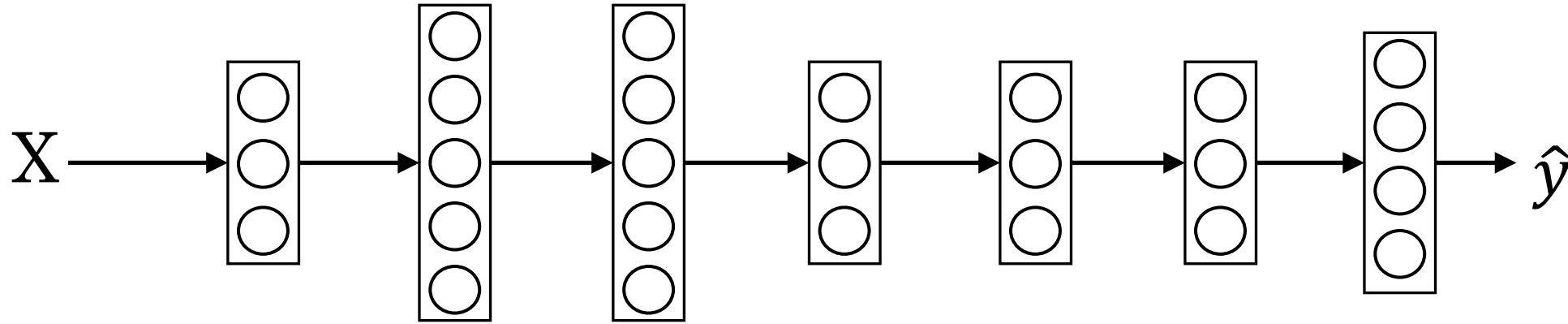
0



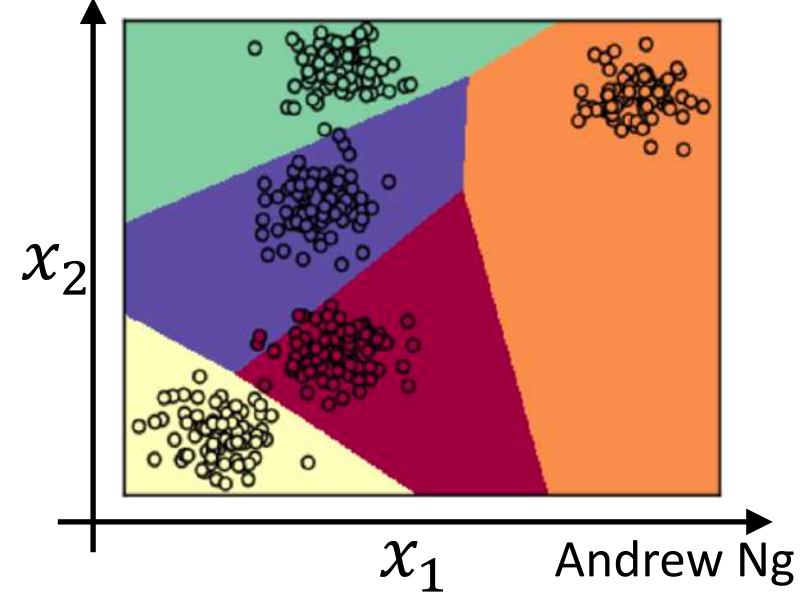
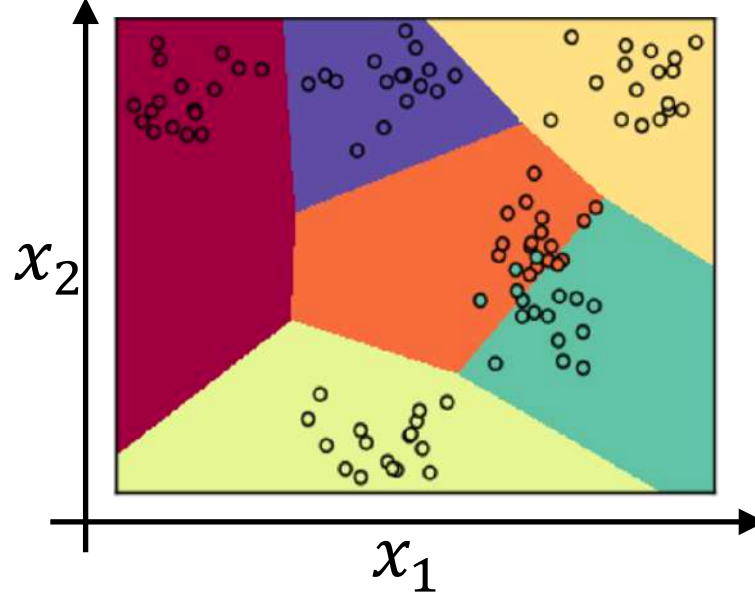
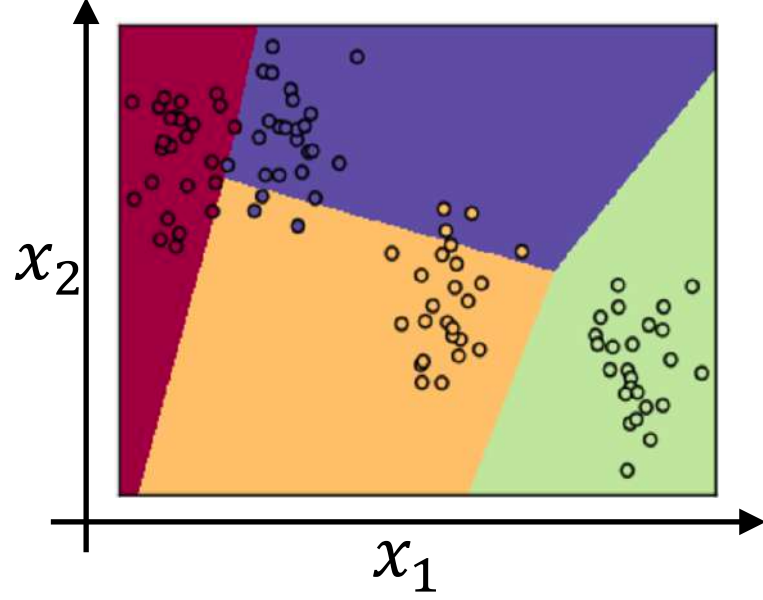
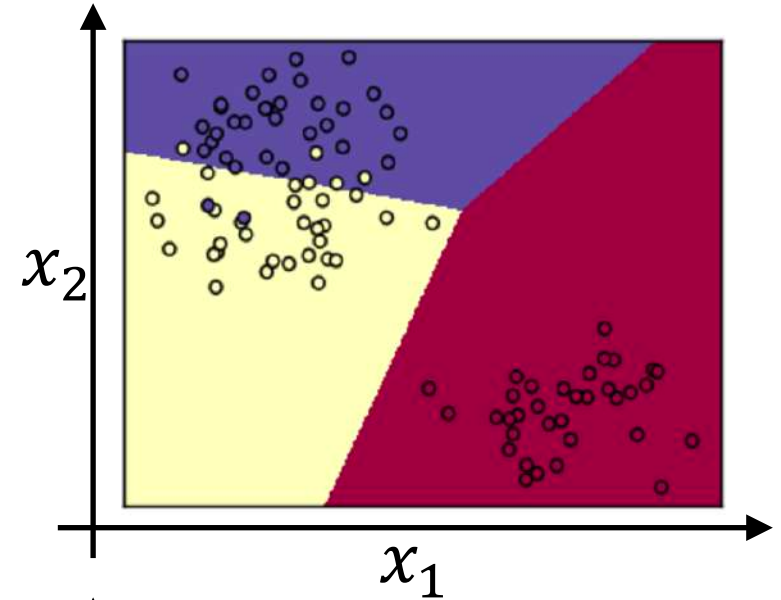
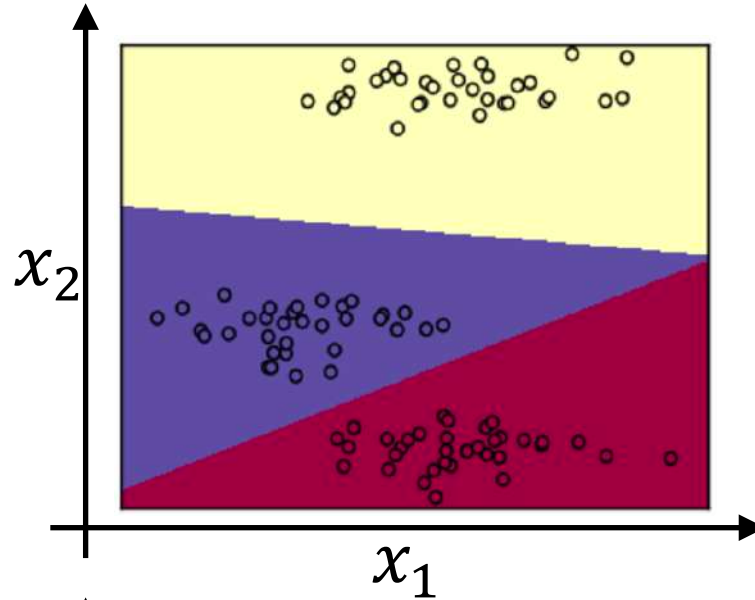
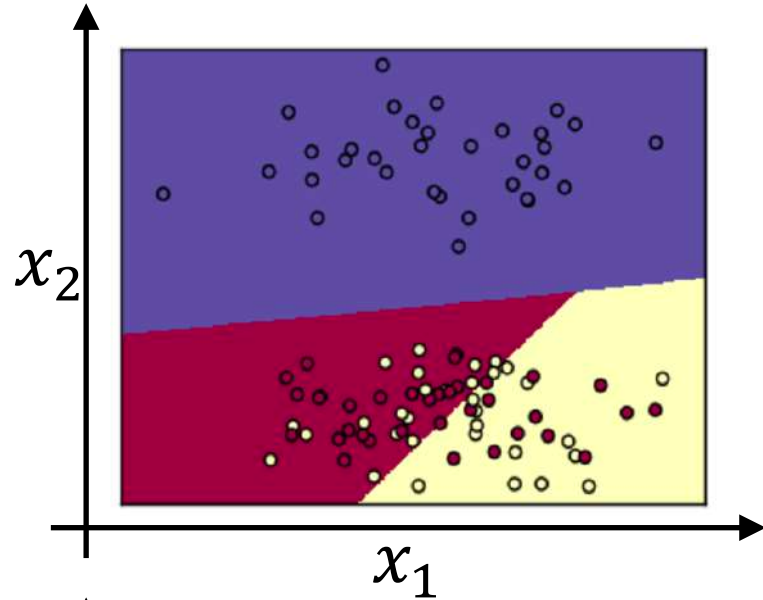
1



Softmax layer



Softmax examples





deeplearning.ai

Programming Frameworks

Deep Learning frameworks

Deep learning frameworks

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

Choosing deep learning frameworks

- Ease of programming (development and deployment)
- Running speed
- - Truly open (open source with good governance)



deeplearning.ai

Programming Frameworks

TensorFlow

Motivating problem

$$\begin{aligned} J(w) &= \boxed{w^2 - 10w + 25} \\ &\quad \swarrow \\ &\quad (w-5)^2 \\ &\quad w=5 \end{aligned}$$

$$\begin{aligned} J(w, b) \\ \uparrow \quad \uparrow \end{aligned}$$