

Hyper Parameters Tuning:

Importance of Hyper - parameters:

High α : learning rate

importance β : momentum term (0-1)

hidden units

mini-batch size

layers

↓ low γ : learning rate decay

$\beta_1, \beta_2, \epsilon$: momentum in the Adam optimization

The way to choose Hyper-parameters:

① In Machine Learning, when there are only a few hps, we use grid-search to systematically choose their best value.

Hp 2

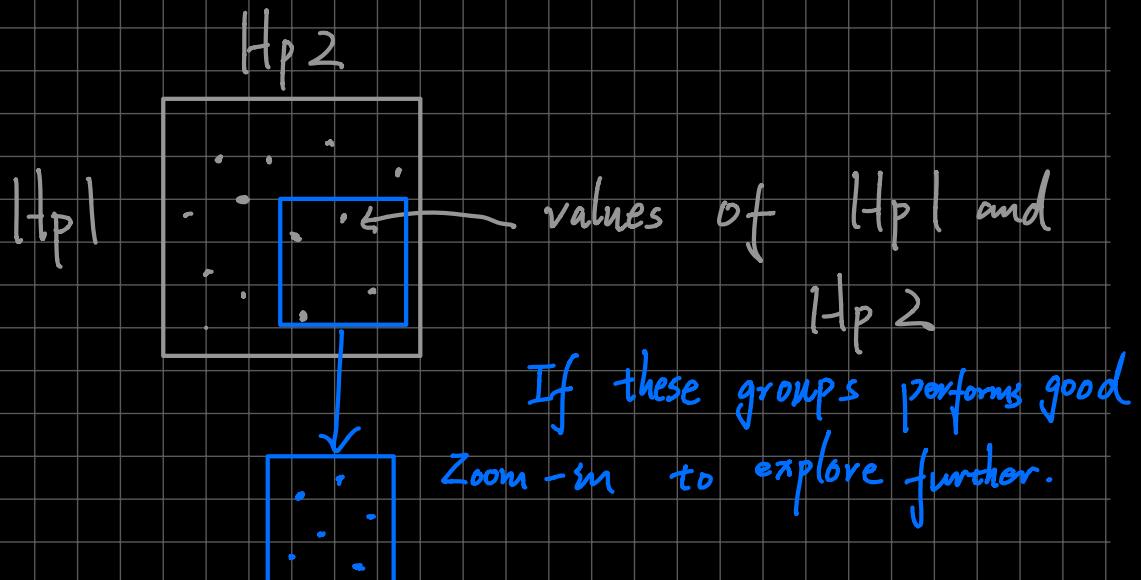
•	•	•	•	•
•	•	•	•	•
•	•	•	•	•
•	•	•	•	•
•	•	•	•	•

values of Hp 1 and

Hp 2

② In DeepLearning, there are a great number of hps and each group of hps need to be trained for a long time. So we don't use this method.

We use Random-Search method to choose -

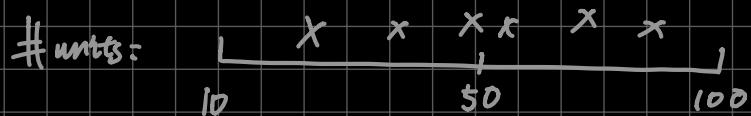


Reason: If use the grid-search method, 25 groups of hps only have 5 different Hp1 values, while in Random-search we have 25 different Hp1 values.

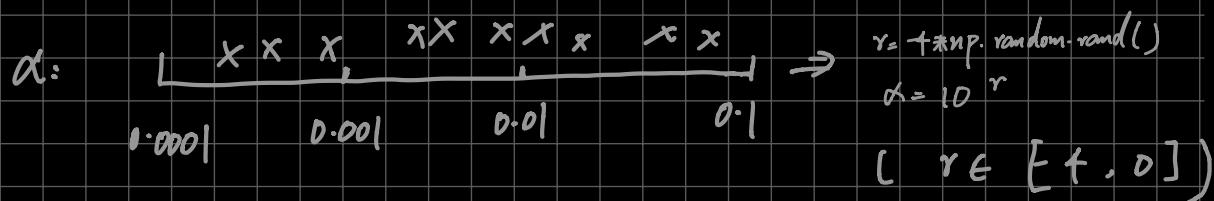
Using an appropriate scale \Rightarrow

Random-choose doesn't mean
that you choose points uniformly
at random.

Linear Scale: # layers . # units



Log Scale: α



Exponentially weighted averages = β



$$\Rightarrow \underbrace{1-\beta}_{\sim} = 0.1, 0.01, 0.001$$

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

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

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

△ why we don't use linear scale for β :

$$\beta = \underbrace{0.9000}_{\sim 10} \rightarrow \underbrace{0.9005}_{\sim 10} \Rightarrow \text{Have no impact}$$

Last Examples: ~ 10

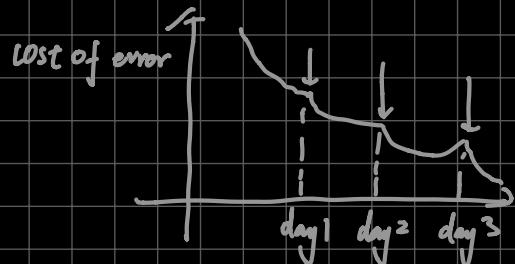
$$\frac{1}{1-p}$$

$$\beta = \underbrace{0.999}_{\sim 1000} \rightarrow \underbrace{0.9995}_{\sim 2000} \Rightarrow$$

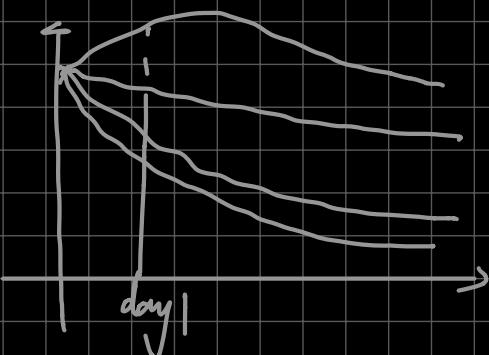
Last Examples = ~ 1000

[Hyper-parameters tuning] method:

pandas: Babysitting one model



canvav: Training many models in parallel.



Implementing Batch Normalization:

Given some intermediate values: $z_1^{(k)}, z_2^{(k)}, \dots, z_m^{(k)}$,

$$\Rightarrow \textcircled{1} u = \frac{1}{m} \sum_{i=1}^m z_i^{(k)} \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad \left. \begin{array}{l} z_{\text{norm}}^{(i)} = \frac{z_i^{(k)} - u}{\sqrt{\sigma^2 + \epsilon}} \\ \sigma^2 = \frac{1}{m} \sum_{i=1}^m (z_i^{(k)} - u)^2 \end{array} \right\}$$

$$\textcircled{2} \quad \tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta \quad \text{learnable parameters}$$

—— Normalization

—— Linear Transformation \Rightarrow 稳定-损失的
层参数信息。

Batch norm is usually used accompanied by mini-batch.

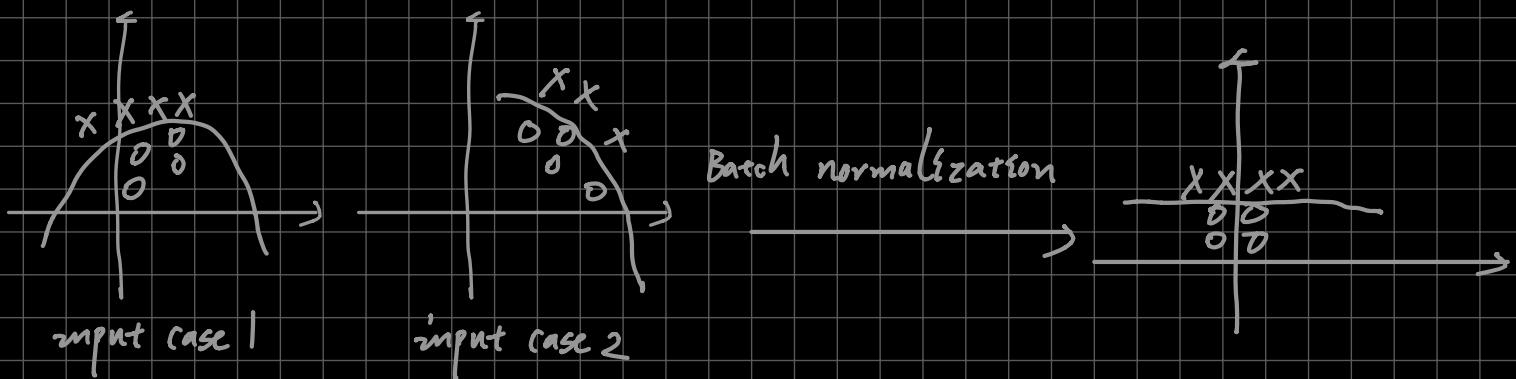
Parameters during training: $W \cdot \underbrace{b}_{\text{over}} \cdot \underbrace{\beta}_{\text{instead}} \cdot \gamma$

$$\beta^{(l)} = \beta^{(l)} - \alpha \cdot d\beta^{(l)}$$

$$\gamma^{(l)} = \gamma^{(l)} - \alpha \cdot d\gamma^{(l)}$$

Why Batch-Normalization works?

- Shift different input distribution to the same norm distribution. Increase the generalized ability of the model (\Rightarrow)



- Avoid ICS (Interval Covariate Shift)

\Rightarrow By ensuring the $u^{(k)} \cdot z^{(k)}$ the same no matter how the $a^{(k-1)}$ changes. (The earlier layer changes)
Allow each layer to learn by itself, reduce the impact of earlier layers.

- Avoid gradient saturation. \Rightarrow Similar to relu / dropout, Batch-normalization adds some noise to $z^{[k]}$ within that minibatch. So it has a slight regularization effect.

Δ mini-batch size \uparrow , regularization effect \downarrow

$64 \rightarrow 512$, we may need to add drop out.

(because the noise we add into $Z^{[L]}$ reduces)

So don't expect to use Batch Normalization as a regularizer. The regularization effect of BN is unintended side-effect.

Batch norm at test time \Rightarrow In training step, we have mini-batch, which Batch norm relies on. while in the test step - we usually feed in the whole test set. Thus it need to do some changes.

Idea: Estimate μ, σ^2 from training set.

\Rightarrow Use exponentially weighted average across

mini-batch as μ, σ^2 estimation value.

Method: During training, keep track of each layer's $\mu^{(l)}, \sigma^2(l)$ and apply to test set in:

$$Z_{\text{norm}} = \frac{z - \hat{\mu}}{\sqrt{\hat{\sigma}^2 + \epsilon}} \text{ when tests.}$$

This method is quite robust.

Multi-class classification problem \Rightarrow Use softmax regression

Binary classification problem \Rightarrow Logistic regression

Regression problem \Rightarrow Linear regression

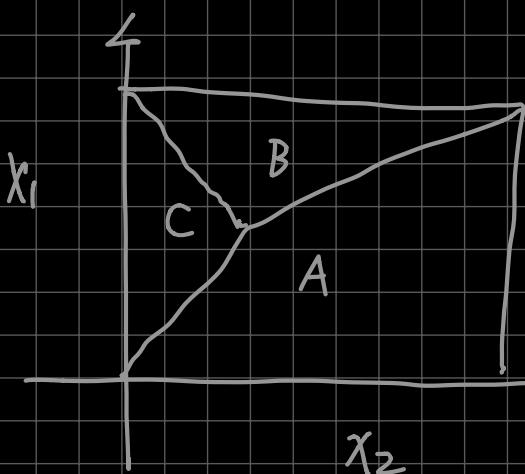
Multi-variate regression.

Softmax function:

$$a^{[L]} = g^{[L]}(z^{[L]})$$

$$t = e^{(z^{[L]})}, t_i = e^{(z_i^{[L]})}$$

$$a^{[L]} = \frac{t}{\sum_{j=1}^n t_j}, a_i^{[L]} = \frac{t_i}{\sum_{j=1}^n t_j}$$



\Rightarrow After using softmax function,
each class is separated by
straight lines.

Why it called softmax ?

Contrast to hard-max (one-hot). where
there is : $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ in the output.

Softmax Regression:

Loss function: $L(y, \hat{y}) = -\sum_{j=1}^n y_j \log \hat{y}_j \Rightarrow \hat{y}_j \rightarrow 1, y_j \rightarrow 1$
or $\hat{y}_j \rightarrow 0, y_j \rightarrow 0$
使得 $L \rightarrow 0$.

Cost function: $J(w^{[i]}, b^{[i]})$

$$= \frac{1}{m} \sum_{i=1}^m L(\underbrace{y_i}_{n \times 1}, \underbrace{\hat{y}_i}_{n \times 1})$$

$$Y = \underbrace{\downarrow \left[\begin{array}{c} \cdots \\ \vdots \\ \cdots \end{array} \right]}_n \xrightarrow{n \times m} n$$

$$\hat{Y} = \underbrace{\downarrow \left[\begin{array}{c} \cdots \\ \vdots \\ \cdots \end{array} \right]}_n \xrightarrow{n \times m} n$$