

Week 3a - Hyperparameter tuning

Friday, August 14, 2020 11:29 PM

① Tuning Process

→ Hyperparameters

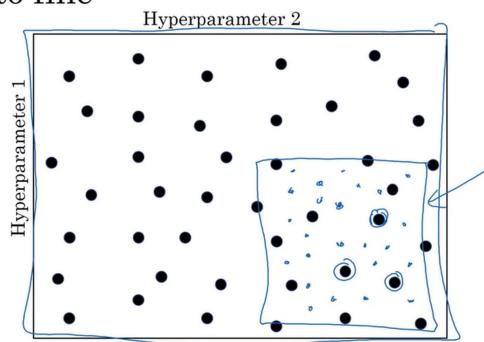
- learning rate α
- $\beta_1, \beta_2, \epsilon^{10^{-8}}$
- no. of layers
- no. of hidden units
- learning rate decay
- mini-batch size

■ → most important to tune
■ → moderate importance
■ → least important

→ Try random parameters : Don't use a grid

→ Use coarse to fine sampling

Coarse to fine



② Using an appropriate scale to pick hyperparameters

→ Picking hyperparameters at random

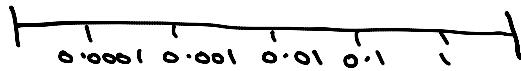
- For layers we might search 1, ..., 4
- for $n^{10^3} \rightarrow 50, \dots, 100$ might be reasonable

→ appropriate scale for hyperparameters

- $\alpha = 0.0001, \dots, 1$
 ↓ small ↑ large

- searching at random doesn't work since 0% values are between 0.1 & 1
- only 10% between 0.0001 & 1

- Hence we use logarithmic scale



• Implementing in python

$$\begin{aligned} r &= -4 * np.random.rand() \rightarrow r \in [-4, 0] \\ \alpha &= 10^r \rightarrow 10^{-4} \dots 10^0 \end{aligned}$$

Generally,

to generate values from 10^a to 10^b
in logarithmic..

generally,

to generate values from 10^a to 10^b
logarithmically

generate 10^r randomly between a & b
and 10^r gives req. values

→ Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \dots 0.999$$

instead use

$$1 - \beta = 0.1 \dots 0.001 \quad 10^{-1} \text{ to } 10^{-3}$$

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

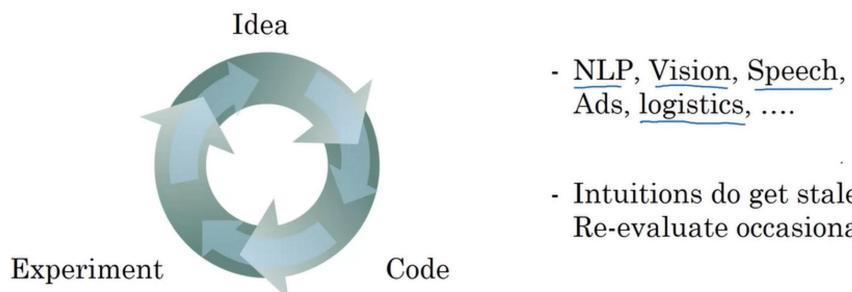
$$1 - \beta = 10^r$$

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

This is done because β is sensitive to small changes when β is larger
(e.g. $0.999 \rightarrow 0.9995$ make more diff than $0.9000 \rightarrow 0.9005$ even though difference is same)

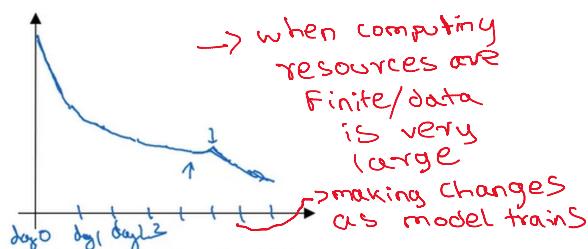
③ Hyperparameters tuning in practice:

Re-test hyperparameters occasionally



Andrew Ng

Babysitting one model



Training many models in parallel

