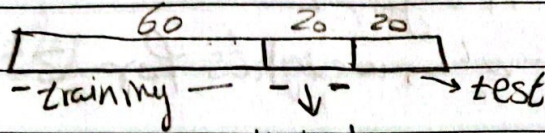


# \* improving DNN, (Hyperparameters, tuning, reg. & optimization) (week 1)

## → Train/Dev/test sets

- when training NN, we have a lot of decisions ex: num of layers
- num of hidden units • learning rate • act. func.

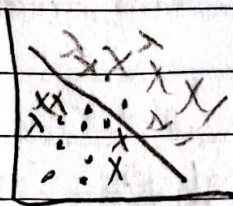


- Holdout  
Cross validation

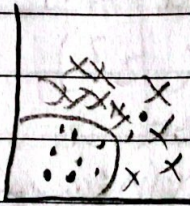
- Development set  
"Dev"

• first training algorithm on training set then use dev to see which of many diff models performs best on your dev set & then after having final model we evaluate the test set

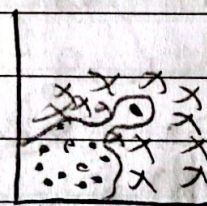
## → Bias / Variance



- high bias
- underfitting



- Just right



- high variance
- overfitting

• train set: 1%

15%

15% → high bias

• Dev set: 11%

16%

30% → high variance

high variance

high bias

↳ when algorithm doesn't do well on the training set.

o. 5% → low bias

1%

→ low variance

} → we search for this



## Basic recipe for ML

- High bias?  
(training data performance)

try

Bigger Neffwort

- train longer
- NN architecture search

after fitting

High Variance?

to evaluate that I look at (dev set performance)

that I have  
good training & dev sets

more data (to train)

try regularization

NN arch. search

- Bias variance trade off: reduce or inc. bias or variance without affecting the other one

## regularization

for ex: in logistic regression we try to minimize the cost fun  $J$ .

$$J(m, b) = \frac{1}{m} \sum \mathcal{L}(\hat{y}_i, y_i) + \frac{\lambda}{2m} \|w\|_2^2$$

regularization parameter

$$\frac{1}{2} \|w\|_2^2 = \sum w_j^2 = w^T w$$

$$\frac{\lambda}{2m} \sum |w_j| = \frac{\lambda}{2m} \|w\|_1$$



→ Drop out regularization:

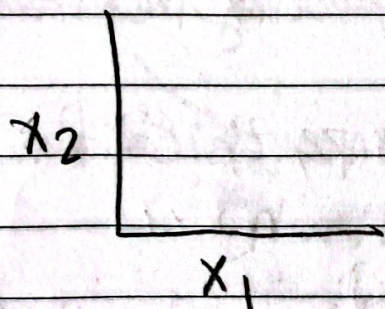
→ Of the regularization techs.

• Early stopping & Data augmentation

→ Normalizing inputs.

• Normalizing the inputs is one of the techniques that is used to speed up the training.

• Normalizing inputs corresponds to



$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \text{ two steps}$$

① to subtract out the mean

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$X = X - \mu$$

② Normalize the Variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} \times x^{(i)} \rightarrow X / \sigma$$

• when features are on similar scales not from 1-1000 & 0-1 but mostly from -1-1 or about similar variance that just makes your cost function easier & faster to optimize

@X:  $x_1: 0 \sim 1$   
 $x_2: -1 \sim 1$   
 $x_3: 1 \sim 2$

} → similar ranges

$x_1: 1 \sim 1000$   
 $x_2: 0 \sim 1$  } hurts optimization



→ Vanishing / exploding gradients  
is one of the problems of training neural network & it is when you're training very deep NN your derivatives or slope can sometimes get either very big or very small & this makes training difficult. & the random weight initialization to significantly reduce this problem.

→ weight initialization for deep networks

$$z = w_1 x_1 + \dots + w_n x_n + b$$

The larger  $n \rightarrow$  smaller  $w_i$

The problem of vanishing or exploding gradients as well as choosing reasonable scaling for how you initialize the weights that makes weights doesn't explode quickly or not decay to zero quickly

→ numerical approximation of gradients

in order to make sure you've got all the details right & internal backprop



→ gradient checking.

is a technique that helps to find bugs in implementation of back propagations many times