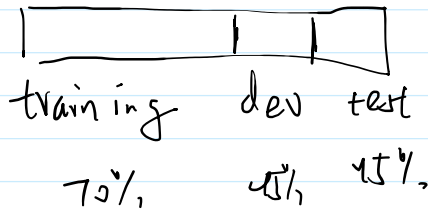


# Weak 1 - Setting of ML Application

Samstag, 18. Juli 2020 14:10

## Training/dev/test sets



Dev: to design/tune hyperparameters

When dataset large enough -> 10,000: 98%, 1%, 1%

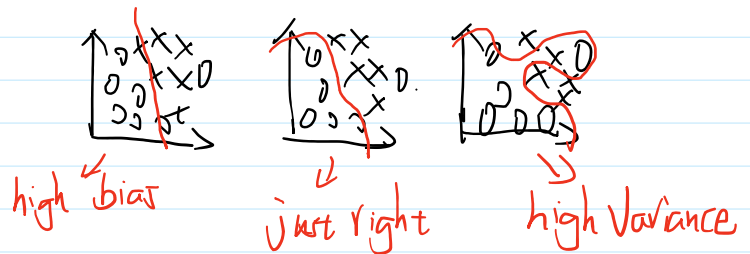
## Evaluation of Algorithms

Bias and variance

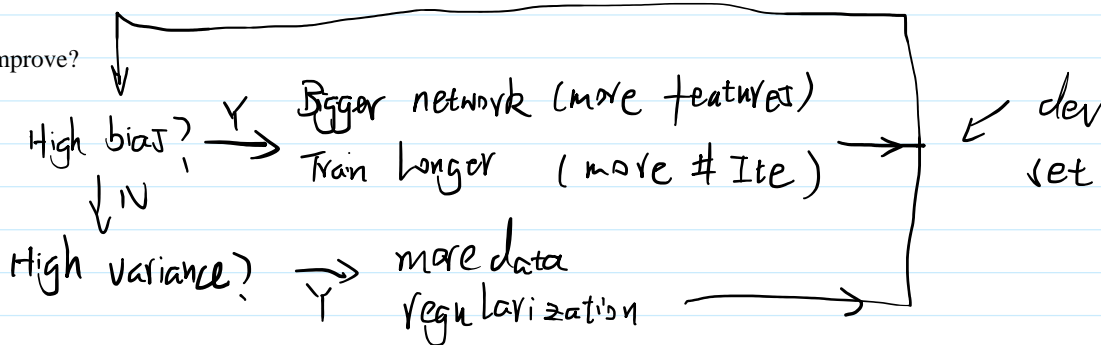
	High variance	High bias	both
Training set error:	1%	15%	15%
Dev set error	11%	16%	30%

High bias: under fitting

High variance: overfitting



How to improve?



Among them, bigger network almost don't hurt variance and more data almost don't hurt bias.

## Regularization

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{y}_i, y_i) + \frac{\lambda}{2m} \|w\|^2 + \left( \frac{\lambda}{2m} b^2 \right)$$

can be omitted

$$\|w\|^2 = \sum_{i=1}^{n_x} w_i^2 = W^T W$$

L2 Regularization

L1 Regularization

$$\underbrace{\sum_{l=1}^L \text{np\_sum}(\text{np\_square}(W^{[l]}))}_{\text{for every layer } L}$$

$$\frac{\lambda}{2m} \|w\|$$

# $\lambda$ regularization parameter

For NN:

$$J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{i=1}^L \|w^{[i]}\|^2$$

$$\|w^{[i]}\|^2 = \sum_{i=1}^{n^{[i-1]}} \sum_{j=1}^{n^{[i]}} (w_{ij}^{[i]})^2$$

$$w^{[L]} = w^{[L]} - \alpha \frac{d}{dw^{[L]}} J = w^{[L]} - \alpha \left[ \underbrace{\left( \text{from back propa} \right)}_{\text{decrease } w} + \underbrace{\frac{\lambda}{2m} w^{[L]}}_{\text{weight decay}} \right]$$

Why regu reduces overfitting?

$$J(\lambda) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{i=1}^L \|w^{[i]}\|^2$$

When lambda large, the goal of minimizing  $J(\lambda)$  is becoming to minimize  $\frac{\lambda}{2m} \sum_{i=1}^L \|w^{[i]}\|^2 \Rightarrow w^{[i]} \approx 0$

Variance  $\downarrow$   $\leftarrow$  Less features  $\leftarrow$  A lot of hidden layer = 0

## Dropout

Keep-prob = 0.8  $\rightarrow$  0.2 dropout

For example layer 3:

`d3 = np.random.randn(a3.shape[0], a3.shape[1]) < keep-prob`

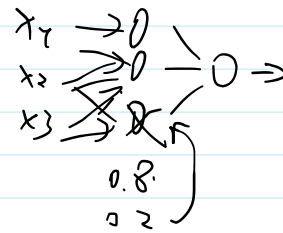
`a3 = a3 * d3` Some units will be 0

`A3 /= keep-prob` # to keep the `a3` and then `z4` in the same scale dont decrease.

Dropout only used for training set, not test set because it will bring in random result.

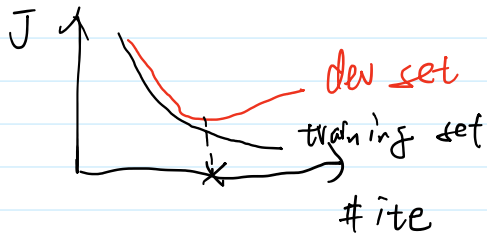
Different dropout (keep-prob) can be used on different layers. If one layer has many units, keep-prob can be smaller.

For computer vision application, dropout is usually default because the input size is so big (so many features,  $X_n$  large), so the training example  $m$  are always relative too small  $\rightarrow$  high variance.

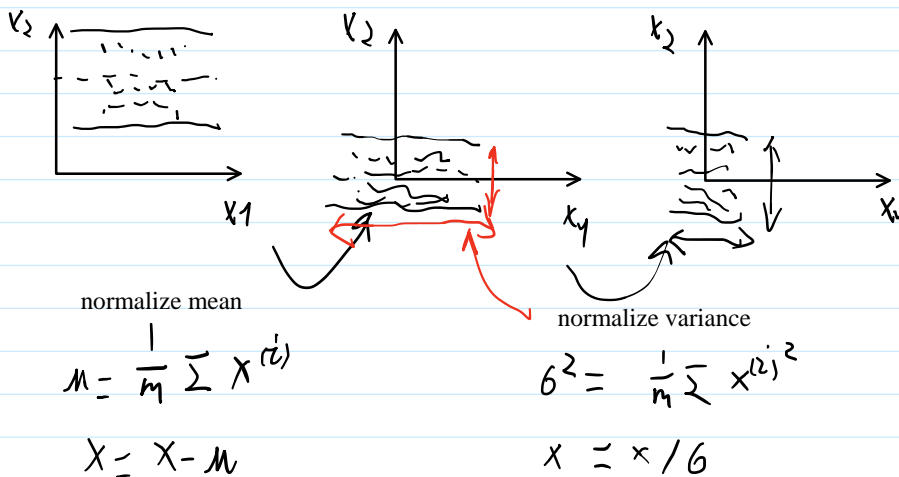


## Other Regularization techniques

1. More input data -> data augmentation
  - a. Horizontal mirror images
  - b. Random rotations
  - c. Strong distortion
2. Early stop



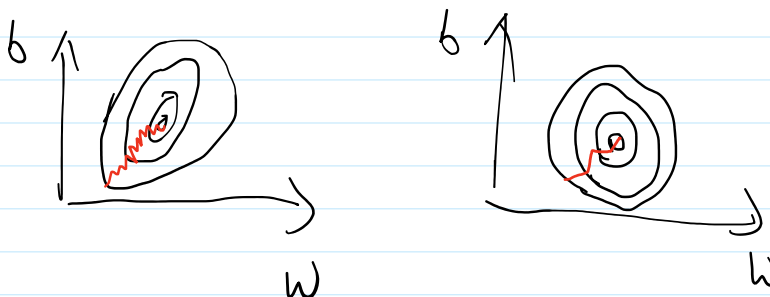
## Data Normalization



use same parameters to normalize train and test set

## When Normalization required?

when data on similar scale, learning is faster, which allows a bigger learning rate.



## Vanishing / Exploding gradients

when  $w^{[l]} > 1$  or  $w^{[l]} < 1$  (identical matrix), activation and gradients will increase or decrease exponentially, which makes learning hard.

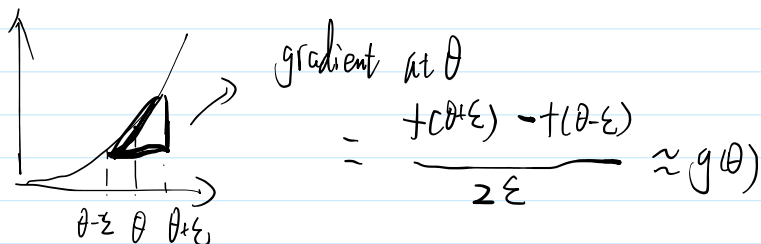
To solve Vanishing / Exploding gradients problem:

Random initialize weight  $w$  for layer  $l$  by:

$$w^{[l]} = \underbrace{\text{np.random.randn}(\text{shape}) * \text{np.sqrt}\left(\frac{2}{n^{[l-1]}}\right)}_{\text{if activation = ReLU}}$$

$$\underbrace{\sqrt{\frac{1}{n^{[l-1]}}}}_{\text{if activation = tanh}}$$

Gradient Checking -- using mathematical relationship to calculate gradient



Implementation:

1. Reshape all parameters ( $w_1, b_1, w_2, b_2, \dots$ ) into vectors and concatenating  $\rightarrow J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = J(\theta)$  them into one vector.
2. Do the previous to all  $dw$  and  $db$  to get  $d_{\theta}$
3. For each item from  $\theta$

$$d\theta_{\text{approx}}^{[i]} = \frac{J(\theta_1, \theta_2, \dots, \theta_{i+\epsilon}, \dots) + J(\theta_1, \theta_2, \dots, \theta_{i-\epsilon}, \dots)}{2\epsilon}$$

should  $\approx d\theta^{[i]} = \frac{\partial J}{\partial \theta^i}$

check  $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx 10^{-7}$  by  $\epsilon = 10^{-7} \Rightarrow \text{great!}$

$10^{-3} \Rightarrow \text{worry!}$

- Don't use gradient check in training, only to debug model.
- Gradient check can't be used with drop off

Take away from the programming assignment:

- Initialization correctly helps the NN model to converge faster and get a lower training error.
  - Initialize w and b as zeros for each layer leads to symmetry fail, meaning the cost will not decrease and model will predict 0 for each examples, because each layer will learn the same thing. So at last it is just a linear classifier such as logistic regression.
  - Set the weights randomly can break the symmetry. small random values performs better than larger.
  - He initialization / Xavier initialization: uses a scaling factor for the weights  $\sqrt{1/\text{layers\_dims}[l-1]}$  (Xavier) ; He initialization would use  $\sqrt{2/\text{layers\_dims}[l-1]}$  and He initialization works well for networks with ReLU activations. -> help to solve vanishing / exploding gradients problem.

```

Def initialize(layers_dims):
    parameters = {}
    L = len(layers_dims) - 1 # integer representing the number of layers

    for l in range(1, L + 1):
        ### START CODE HERE ### (~ 2 lines of code)
        parameters['W' + str(l)] = np.random.randn(layers_dims[l], layers_dims[l-1]) * np.sqrt(2./layers_dims[l-1])
        parameters['b' + str(l)] = np.zeros(shape=(layers_dims[l], 1))
        ### END CODE HERE ###

    return parameters

```

- Regularization: L2 rege or dropout
  - L2 regularization: apply the lambda term on cost function and gradients (dW1, dW2, dW3.....)

```

def compute_cost_with_regularization(A3, Y, parameters, lambda):
    """
    Arguments:
    A3 -- post-activation, output of forward propagation, of shape (output size, number of examples)
    Y -- "true" labels vector, of shape (output size, number of examples)
    parameters -- dictionary containing parameters of the model

    Returns:
    cost
    """
    m = Y.shape[1]
    W1 = parameters["W1"]
    W2 = parameters["W2"]
    W3 = parameters["W3"]

    cross_entropy_cost = compute_cost(A3, Y) # This gives you the cross-entropy part of the cost
    L2_regularization_cost = 1/m * lambda/2 * (np.sum(np.square(W1)) + np.sum(np.square(W2)) + np.sum(np.square(W3)))
    cost = cross_entropy_cost + L2_regularization_cost

    return cost

```

- dropout
  - in forward propagation
 

```

D1 = np.random.rand(A1.shape[0], A1.shape[1]) # Step 1: initialize matrix D1 = np.random.rand(..., ...)
D1 = (D1 < keep_prob).astype(int)              # Step 2: convert entries of D1 to 0 or 1 (using keep_prob as threshold)
A1 = A1 * D1                                   # Step 3: shut down some neurons of A1
A1 = A1 / keep_prob                             # Step 4: scale the value of neurons that haven't been shut down
          
```
  - in backward:
 

```

          shut down the same neurons, by applying the same mask D[1] to dA1.
          divide dA1 by keep_prob again
          
```
  - only apply drop out on training set because we don't want randomness in the prediction

- Gradient Check
  - mathematically calculate gradient and then compare it with the result from the backpropagation.
  - Before applying gradient check you have to be sure that the J cost function is computed in a right way, because the gradient check based on that.
  - derivative:

$$\frac{\partial J}{\partial \theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

For each  $i$  in `num_parameters`:

- To compute `J_plus[i]`:
  1. Set  $\theta^+$  to `np.copy(parameters_values)`
  2. Set  $\theta_i^+$  to  $\theta_i^+ + \epsilon$
  3. Calculate  $J_i^+$  using `forward_propagation_n(x, y, vector_to_dictionary( $\theta^+$ ))`.
- To compute `J_minus[i]`: do the same thing with  $\theta^-$
- Compute  $gradapprox[i] = \frac{J_i^+ - J_i^-}{2\epsilon}$

Thus, you get a vector `gradapprox`, where `gradapprox[i]` is an approximation of the gradient with respect to `parameter_values[i]`. You can now compare this `gradapprox` vector to the gradients vector from backpropagation. Just like for the 1D case (Steps 1', 2', 3'), compute:

$$difference = \frac{\|grad - gradapprox\|_2}{\|grad\|_2 + \|gradapprox\|_2}$$