

Week 1b - Regularizing your Neural network

Friday, August 14, 2020 11:25 PM

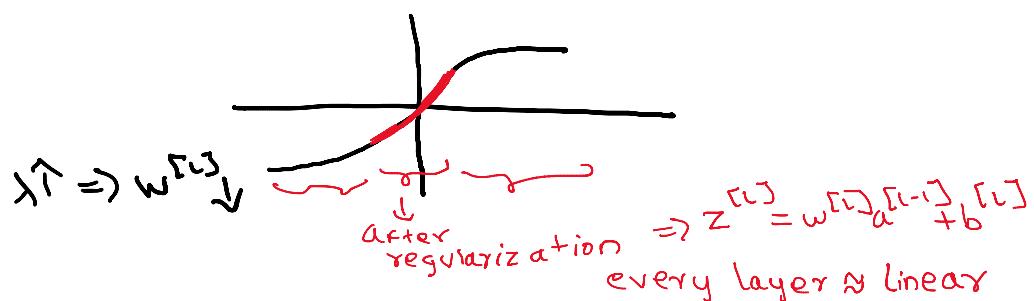
① Regularization

- For logistic regression $\min_{w,b} J(w,b)$ $\text{lambd} \leftarrow \lambda \rightarrow \text{regularization parameter}$
 $J(w,b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2 + \frac{\lambda}{2m} b^2$
L₂ regularization: $\|w\|_2^2 = \sum_{i=1}^{n_{\text{vec}}} w_i^2 = w^T w \rightarrow \text{most used}$
L₁ regularization: $\frac{\lambda}{m} \sum_{i=1}^m |w_i| = \frac{\lambda}{2m} \|w\|_1 \rightarrow w \text{ will be sparse (lot of zeros)}$
- For a neural network
 - $\Rightarrow J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{[i]}, y^{[i]}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$
where, $\|w^{[l]}\|_F^2 = \sum_{i=1}^{n_l} \sum_{j=1}^{n_{l-1}} (w_{ij}^{[l]})^2$ $w: [n^{[0]}, n^{[L-1]}]$
"Frobenius norm"
 - $\Rightarrow dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}$ $\frac{\partial J}{\partial w} = dw^{[l]}$
 $\Rightarrow w^{[l]} := w^{[l]} - \alpha dw^{[l]}$
"weight decay" $w^{[l]} := w^{[l]} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \right]$
 $= w^{[l]} - \underbrace{\alpha \frac{\lambda}{m} w^{[l]}}_{w^{[l]}(1-\frac{\alpha \lambda}{m})} - \alpha (\text{from backprop})$
 $\Rightarrow w^{[l]}(1-\frac{\alpha \lambda}{m}) > \text{slightly less than } 1$

② Why does regularization reduce overfitting

\rightarrow Intuitively as λ gets bigger, many $w^{[i]}$ get closer to zero and cancel their effect leading to simpler network

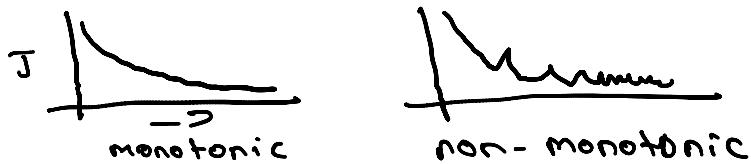
\rightarrow In tanh function,



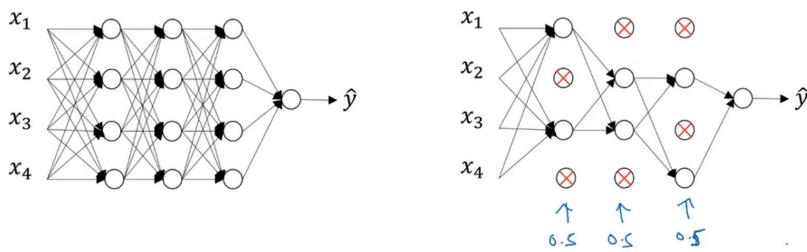
$\rightarrow \lambda$ can be considered as a penalty

every layer is linear

- λ can be considered as a penalty for the parameter being too large
- It makes J reduce monotonically



③ Dropout regularization



eliminate nodes at random to reduce size of network \Leftrightarrow thereby reducing overfitting

- Implementing dropout ("Inverted dropout")

Consider a random layer (3 in this case)

dropout vector $d_3 = np.random.rand(a3.shape[0], a3.shape[1]) < \text{keep_prob}$

$a_3' = np.multiply(a3, d_3)$ → zeroes out

$a_3' /= \text{keep_prob}$ nodes where $d_3 \leq 0$

prob. of elimination

- d_3 will be a boolean array
e.g.: if $\text{keep_prob} = 0.8$,
the prob. of a node not being eliminated is 80%
 $\Rightarrow 50$ units $\Rightarrow 40$ units

value of z_{4j}' is lesser
hence we divide by keep_prob
to adjust for it

- Making predictions at test-time

$$a^{[0]} = x$$

- dropout isn't used in test as it has already influenced training

- If dropout is used in prediction, different outputs are obtained every time

④ Understanding dropout

④ Understanding dropout

- Intuition:
- Can't rely on any one feature, so have to spread out weights
 - Since any feature can disappear, the neuron doesn't decide too much based on one input
 - Results in shrinking weights and spreaded out weights

Note:

- $\text{keep_prob} = 1$ generally for input since we don't want to lose useful information

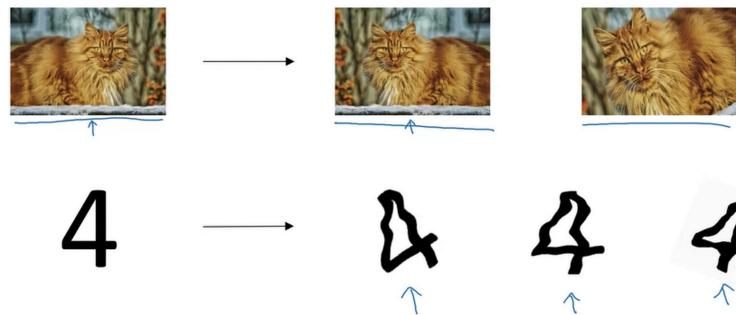
- Used only if overfitting is a problem
eg: in computer vision

→ Cost function is no longer well-defined since every time a different function is computed

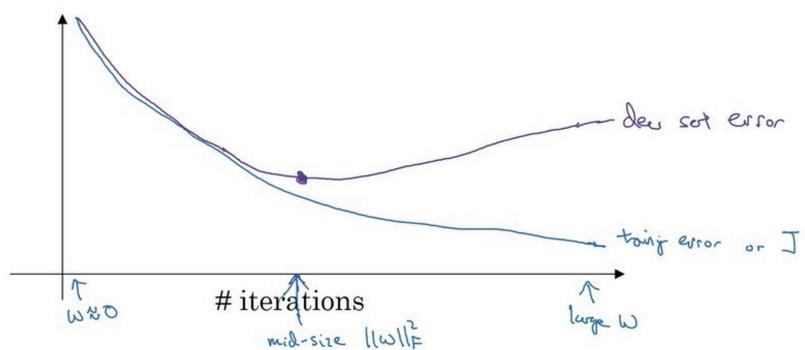
⑤ Other regularization methods

→ Data augmentation

- By flipping images, we can effectively double the no. of data
- Random distortions / zooms can also be applied.



→ Early stopping



- ML has 2 tasks primarily, optimization + regularization

- ML has 2 tasks primarily,
known as Orthogonalization
 - i. Optimize cost fn.
- gradient ... }
ii. Not overfit
- regularization ... }

Early stopping combines these two
and limits our options and
complicates the tuning process