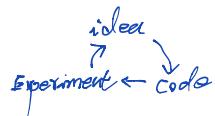


- Setting up your ML application

- Train / dev / test sets

Applied ML is a highly iterative process

layers, # hidden units, learning rates, activation functions, ...



Not having a test set might be okay (only dev set (cv))

- Bias / Variance

- Basic "recipic" for machine learning

High bias? →

- Bigger Network

- Train longer

High Variance? →

- More data

- Regularization

- NN architecture search

- Regularizing your Neural Network

- Regularization

o Logistic Regression:

L2 Regularization

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, g^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

L1 Regularization

$$J(w, b) = \dots + \frac{\lambda}{2m} \|w\|_1$$

w will be sparse (a few 0s in w)

λ : Regularization Parameter

o Neural Network

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, g^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} (w^{[l]}_{ij})^2 \quad \text{--- Frobinius Norm of Matrix}$$

$$dw^{[l]} = (\text{from back prop}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l]} := w^{[l]} - \alpha \underbrace{dw^{[l]}}$$

$$= w^{[l]} - \underbrace{\frac{\alpha \lambda}{m} w^{[l]}}_{\text{Weight decay}} - \alpha (\text{from back prop})$$

"Weight decay"

o Why regularization Reduces Overfitting

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, g^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

Reducing the impact of some of the hidden unit,

by adding penalize term. $\frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$

o Dropout Regularization

$$\begin{matrix} x_1 & 0 & 0 & 0 \\ x_2 & 0 & 0 & 0 \\ x_3 & 0 & 0 & 0 \\ x_4 & 0 & 0 & 0 \end{matrix} \xrightarrow{0 \rightarrow \hat{y}}$$

o Inverted Dropout

Illustrate with $L=3$, $\text{keep-prob} = 0.8$

$$d_3 = \text{np.random.rand}(a_3.shape[0], a_3.shape[1]) < \text{keep-prob}$$

$$a_3 = a_3 * d_3$$

$$a_3 = a_3 / \text{keep-prob} \quad (\text{ensure expected value be the same})$$

o Making predictions at test time

No dropout for testing. otherwise output will be noisy (random)

o Understanding Dropout

Intuition: Can't rely on any one feature, so have to spread out weights
→ shrinking weights

keep-prob: different layers can have different keep-prob values, usually
the large neuron layer may have smaller keep-prob.

Dropout is frequently used in Computer Vision

Downside: J is less well defined

- Other Regularization Method.

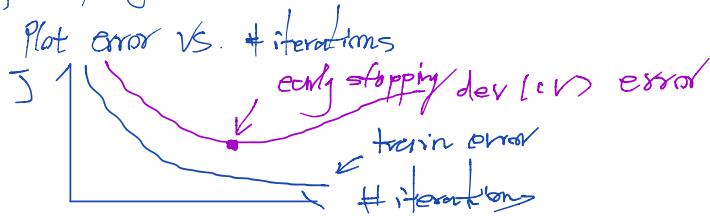
- Data augmentation.

- flip the image to get more training examples

- Random Crops, distortions, rotations of images

- & cheap way to get more data

- Early stopping



Downside: Early stopping couples $\begin{cases} \text{optimize cost func } J \\ \text{- Not overfit} \end{cases}$ orthogonalization

cannot work on these two aspect individually
Make things more complicated

- Setting up your optimization Problem

- Normalizing Inputs

- Normalizing training sets

- Subtract mean

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

- Normalize Variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)})^2 \quad [\text{element-wise, already subtracted mean}]$$

$$x = \frac{x}{\sigma}$$

- Use same μ, σ to normalize test set

- why normalize Inputs

: normalized gives all input features ~~in~~ ^{element-wise} cost function same scale ⁱⁿ ~~out~~ have small learn rate

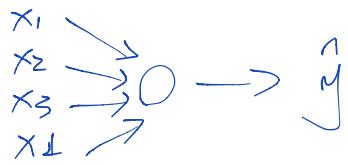
- Vanishing / exploding gradients

An intuitive example:

for a large deep NN, if each $w^{[l]}$ is very long than diag(1), then output approx. $\hat{y} = w^{[L]} \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}^{L-1} x \rightarrow \hat{y}$ explod
if $w^{[l]}$ are all very small than diag(1), $\rightarrow \hat{y}$ vanish
Gradient descent may be very slow

- Weight Initialization for Deep Networks

- Single neuron example



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

large $n \rightarrow$ smaller w_i

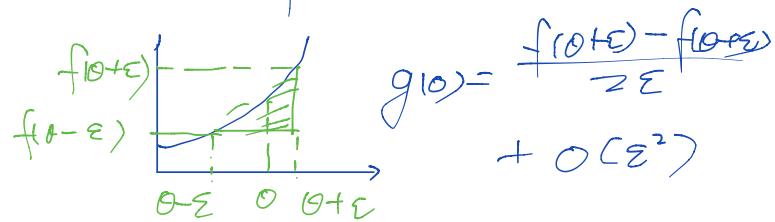
$$\text{Initialize } w_i \text{ so that } \text{Var}(w_i) = \begin{cases} \frac{1}{n} & \text{if tanh} \\ \frac{2}{n} & \text{if ReLU} \end{cases}$$

$$w^{[1]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{\frac{1}{n}}{\frac{2}{n}}\right)$$

(according to published paper)

- Numerical Approximation of gradients

- check your derivative computation



$$\text{Sekante: } f'(0) = \lim_{\epsilon \rightarrow 0} \frac{f(0+\epsilon) - f(0-\epsilon)}{2\epsilon} = \lim_{\epsilon \rightarrow 0} \frac{f(0+\epsilon) - f(0)}{\epsilon}$$

$$\epsilon < 1 \quad f'(0) = \frac{f(0+\epsilon) - f(0-\epsilon)}{2\epsilon} + O(\epsilon^2)$$

$$f'(0) = \frac{f(0+\epsilon) - f(0)}{\epsilon} + O(\epsilon)$$

1st implementation yields smaller estimation error

- Gradient checking
 - Gradient checking for a neural network
 - Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector θ
 - $J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$
 - Take $d\theta^{[1]}, d\theta^{[2]}, \dots, d\theta^{[L]}$ and reshape into a big vector $d\theta$
- Gradient checking
 - for each i :
 - $d\theta_{approx}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$
 - check $d\theta_{approx} \approx d\theta$ (from backprop)
 - $\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2} < 10^{-7}$ reg \Rightarrow ok!
- Grad checking Implementation Notes
 - Don't use in training — only to debug
make sure backprop works as expected
 - If algorithm fails grad check, look at component to identify bug
 - Remember regularization term in $d\theta$ from J

- Grad check doesn't work with Dropout
difficult to calculate ∇ with dropout.
Recommendation: turn off dropout, do Grad check,
make sure algorithm works, then turn on
dropout, hope for the best!
- Run at random initialization; perhaps again after some training