

## Metric

Class prob pred: Log Loss

Regres: RMSE =  $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$  square to ignore pos, neg term

Bias: inability of model to learn from data.

Variance: diff in fits between datasets

Low bias: learning well with dataset

High var: high diff of fitting datasets

Gradient learning rate

$$w_1 = w_0 - \alpha \nabla \frac{1}{2} \sum_{i=1}^n ((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)})^2$$

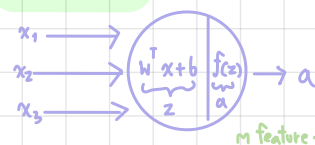
Stochastic learning rate

$$w_1 = w_0 - \alpha \nabla \frac{1}{2} \sum_{i=1}^n ((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)})^2$$

Mini batch learning rate 100 from 1k data

$$w_1 = w_0 - \alpha \nabla \frac{1}{2} \sum_{i=1}^n ((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)})^2$$

neural net



$$z = b + x_1 w_1 + x_2 w_2 + \dots + x_n w_n = b + x^T W$$

Sigmoid:  $f(z) = \frac{1}{1+e^{-z}}$  (prob 0-1 X in class)

Relu:  $\text{relu}(z) = \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases} = \max(0, z)$

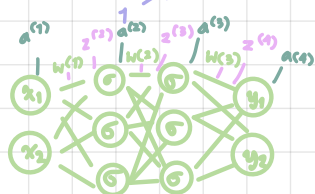
Tanh:  $\tanh(0) = 0, \tanh(\infty) = 1, \tanh(-\infty) = -1$

LReLU:  $\text{LReLU}(z) = \begin{cases} z, & z \geq 0 \\ a, & z < 0 \end{cases} = \max(a, z)$  [a < 1]

(Leaky) LReLU(0) = 0, LReLU(1) = 1, LReLU(-1) = -a

Softmax: multiclass (prob 0-1)

$$z = 1.6 + 0.6 + 0.3 + 0.5 = 2.6$$



$w^{(1)}$ : 2x3 matrix

$z^{(2)}$ : 3 vector

$a^{(2)}$ : 3 vector

## step

Input: vector  $x, z^{(2)} = x W^{(1)}, a^{(2)} = \sigma(z^{(2)})$   
 $z^{(3)} = a^{(2)} W^{(2)}, a^{(3)} = \sigma(z^{(3)}), z^{(4)} = a^{(3)} W^{(3)}$   
 $\hat{y} = \text{softmax}(z^{(4)})$

Loss func: 1 iter, Loss cost: whole iter

Hyper para: upper layer para, tunable

Loss func, Learning rate, epoch, batch size, optimizer

Forward propagate m: min output node (1)  
n: min input node (2)

$W$  size:  $M \times n, X = n \times 1, B = M \times 1$

$$z = W \cdot x + B, z = [w_1, w_2] \cdot [x_1] + [b] = [w_1 x_1 + w_2 x_2 + b]$$

Backpropagation

$$\text{Error at } w_1 = \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

## CNN Steps

- 1) init: weight, bias randomly
- 2) forward: pass input through network layer by layer. compute weighted sum of inputs & apply activation func. Continue to final layer & gen pred output
- 3) Calculate Error: cal error (regression, cross-entropy loss)
- 4) Back: Cal gradient of the loss with respect to weights & biases. update in the opposite direction to minimize loss
- 5) Repeat front/back step for epoch, adjust hyperparams

## Vanishing Gradient

$$\frac{\partial L}{\partial w^{(1)}} = (y - \hat{y}) \cdot \sigma'(z^{(2)}) \cdot w^{(2)} \cdot \sigma'(z^{(3)}) \cdot w^{(3)} \cdot \sigma'(z^{(4)}) \cdot x$$

$$\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z)) \leq 0.25$$

• More layer, smaller gradient on early layers

• Solve w/ ReLU

## Keras

• Sequential Model: linear stack of layers, simpler, more convenient if model've this form

Functional API: more detailed, complex allows more complicated architectures

from keras.models import Sequential  
model = Sequential()

from keras.layers import Dense, Activation

first model.add(Dense(units=4, input\_shape=(x\_train.shape[1:], activation='relu'))

subsequent layer

model.add(Dense(units=4, activation='relu'))

model.add(Dense(units=1, activation='sigmoid'))

## Multiclass

- one hot encoding for categories
- vector's length = categories's num

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

- Final layer's length = pos class's num
- softmax = [0-1] prob
- loss: categorical cross entropy

$$\text{"Log Loss"} \quad CE = - \sum_{i=1}^n y_i \log(\hat{y}_i)$$

$$\frac{\partial CE}{\partial \text{softmax}} = \frac{\partial \text{softmax}}{\partial z_i} = \hat{y}_i - y_i$$

## input scale

Linear Scaling to [0,1]

$$x_i = \frac{x_i - x_{min}}{x_{max} - x_{min}}$$

$$x_i = 2 \left( \frac{x_i - x_{min}}{x_{max} - x_{min}} \right) - 1$$

Standardization var approx. std. normal

$$x_i = \frac{x_i - \bar{x}}{\sigma}, \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

## Regularizing

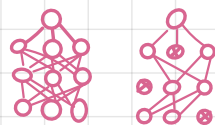
- Penalized cost func

add penalty to loss func for having high weights

$$J = \frac{1}{2n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + \lambda \sum_{j=1}^m w_j^2$$

- Dropout

- randomly remove subset of neurons at each batch
- prevent NN from rely too much on individual path way



training time: present w/ prob p

test time: always present

- Early Stopping

- make a rule to be called to stop the training at some condition.
- Example: check validation log-loss every 10 epochs, if it is higher than it was last time, stop and use previous model

## optimizer

• Momentum: change direction a little each time, keeps a "running average" of step directions, smoothing out the variation of individual points.

$$V_t = \beta \cdot V_{t-1} + (1 - \beta) \cdot \nabla J$$

$$W = W - \alpha \cdot V_t$$

$$(\beta < 1)$$

gradient update'll be bit smoother & more able to travel down shallow directions

- Nesterov Momentum

- Control "overshooting" by looking ahead
- Apply gradient to "non-momentum" component.

$$V_t = \beta \cdot V_{t-1} + \alpha \nabla J(W - \beta \cdot V_{t-1})$$

$$W = W - V_t$$

## Adagrad (Adaptive Grad)

- Scale update for each w separately
- keep running sum of previous updates
- Decide new update with prev sum

$$W = W - \frac{\alpha}{\sqrt{G_t + \epsilon}} \nabla J \quad \text{dynamic learning rate} \quad G_t = G_{t-1} + (\nabla J)^2$$

- learning / sqrt of sum of each component separately

## CNN

- Edges  $\rightarrow$  Shape  $\rightarrow$  relations betw shapes

### kernels (filters)

- grid of weights "overlaid" on image, centered on 1 pix  
traditional image processing techniques

• Blur • Sharpen • Edge detection • Emboss

3 2 1	-1 0 1	
1 2 3	-2 0 2	2
1 1 1	-1 0 1	

$$= (3 \cdot -1) + (2 \cdot 0) + (1 \cdot 1) + (1 \cdot -2) + (2 \cdot 0) + (3 \cdot 2) + (1 \cdot -1) + (1 \cdot 0) + (1 \cdot 1) = 2$$

-1	2	0	3	1
1	0	0	2	2
2	1	2	1	1
0	0	1	0	0
1	2	1	1	1

-1	1	2
1	1	0
-1	-2	0

$$\text{output size} = \left\lfloor \frac{h-f}{s} + 1 \right\rfloor$$

no pad

$$\text{with pad} = \left\lfloor \frac{h+2p-f}{s} + 1 \right\rfloor$$

### Depth (channels)

• RGB: 3 chan, CMYK: 4 chan behavior

• kernel's depth = input's depth

• Computation:  $5 \times 5$  kernel on RGB img  
 $5 \times 5 \times 3 = 75$  weights

output's depth of convo layer =  
num of kernel if applies to prev layer  
each kernel is trained diff

• 10 kernels in layer, layer output depth 10

### local feature detector

-1	1	-1
-1	1	-1
-1	1	-1

-1	-1	-1
1	1	1
-1	-1	-1

-1	1	-1
-1	1	-1
-1	1	-1

- Let NN learn which kernel are most useful
- same set of kernel across entire img (translation invariance)
- Reduce num of params & variance

### Grid size (h & width)

- num of pixels a kernel see at once
- typically odd num so there's "center" pixel

### Padding

valid  $\rightarrow$  no padding  
same  $\rightarrow$  same as input

- prevent edge effect on kernel
- add extra pixel around frame (0)
- every pix of orig img can be center pix as the kernel move

0	0	0	0	0	0
0	1	2	0	3	1
0	-	-	-	-	0
0	-	-	-	-	0
0	-	-	-	-	0
0	-	-	-	-	0
0	-	-	-	-	0
0	0	0	0	0	0

no-pad-out size:  $h-f+1$   
pad: " " :  $h+2p-f+1$

### Pooling

Reduce img size by rapping a patch of pix to single val

• Shrink ing's dimensions

1	1	0	1
1	2	2	5
1	1	3	1
0	1	1	2

Max pooling / Average Pooling

$$\text{pooling layer: } \text{output size} = \frac{n - \text{pool size}}{\text{stride}} + 1$$

## Transfer Learning

• Early layer in NN are hard to train due to "vanishing gradient"

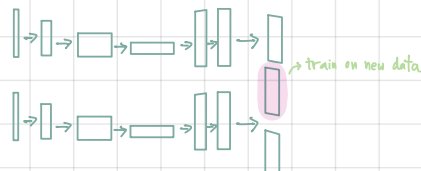
• Feature map: set of conv layer represent specific feature in the input image

• later layer are easier to train with more immediate impact on final result

quicker: back propagating's signal make the earlier weight smaller & smaller

• Famous, competitive model are diff to train from scratch  $\rightarrow$  huge datasets, long num of iteration, very heavy computing machinery, Time experimenting to get hyper-params right

• Idea: keep early layers of pretrained & retrain the later layer for specific UC



Fine tune: train pretrained model for specific use on specific new dataset

### Calculating Conv2d

params: weights + biases

