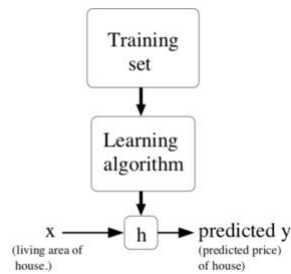


## Cost Function

- Aka Cost function ,error function
- Def: how close the hypothesis functions are to the corresponding y.

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2.$$

- A machine learning hypothesis is a model that approximates the target function and maps the set of inputs to the set of outputs.



## Common Loss functions

- Mean Absolute Error

$$MAE = \frac{1}{m} \sum_{i=1}^m |\hat{y}^{(i)} - y^{(i)}|$$

- Mean Squared Error

$$MSE = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

- Binary cross-entropy  
→ A binary classification

## Gradient descent

- an optimization algorithm used to minimize some function by iteratively moving in the direction of steepest descent as defined by the negative of the gradient
- Y: Cost function(J); X,Z: parameters from hypothesis
- Stochastic gradient descent
  - o update the weights after each training sample

## Activation Function

- a function that you use to get the output of node

$$z = w^T x + b$$

$$a = \sigma(z)$$

## Common Types

- Sigmoid: used when output is between 0 and 1
- ReLu: default
- Tanh: between -1 and 1

## Backward propagation

- output values are compared with the correct answer to compute the value of some predefined error-function

## Bias & variance

- high bias: training set errors are high
  - o Solution
    - Bigger network
    - Advance algo
    - NN architecture search
- high variance: (Dev set errors – training set errors) very large
  - o more data
  - o regularization
  - o NN architecture

## Cross-Validation

- a technique for evaluating ML models by training several ML models on subsets of the available input data and evaluating them on the complementary subset of the data.

## Regularization

- Reduce the test error at the expense of increased training error
- Most regularization strategies are based on regularizing estimators
  - o Works by trading increased bias for reduced variance

## Addressing overfitting

- Manually select which features to keep

## Parameter Norm penalties

- Add a parameter norm penalty to the cost function

$$\tilde{J}(\theta; \mathbf{X}, \mathbf{y}) = J(\theta; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\theta)$$

- $\alpha$  correspond to more regularization
- Works well when there are lots of features
- When our training algorithm minimizes the regularized cost function  $\tilde{J}$ , it will decrease both the original cost  $J$  on the training data and some measure of the size of the parameters  $\theta$  (or some subset of the parameters)
- typically choose to use a parameter norm penalty  $\Omega$  that penalizes only the weights of the affine transformation at each layer and leaves the biases unregularized

## L2 Regularization (more often)

- aka Weight decay
- Drives the weights closer to the origin by adding a regularization term
- Multiplicatively shrink the weight vector by a constant factor on each step

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

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2 \quad w: \begin{pmatrix} n^{[l-1]} & n^{[l-1]} \\ \uparrow & \uparrow \end{pmatrix}$$

- o Forbenius norm

$$\begin{aligned}\underline{w}^{[L]} &:= \underline{w}^{[L]} - \alpha \left[ \left( \text{from backprop} \right) + \frac{\lambda}{n} \underline{w}^{[L]} \right] \\ &= \underline{w}^{[L]} - \underbrace{\left( \frac{\alpha \lambda}{n} \underline{w}^{[L]} \right)}_{\text{L2 regularization}} - \alpha \left( \text{from backprop} \right)\end{aligned}$$

- $\lambda$  : regularization parameters
  - o A hyperparameter
  - o Set this using Dev set
  - o Using hold out cross-validation to find the best values

### L1 Regularization

- The sum of absolute values of the individual parameters

$$\frac{\lambda}{2m} \sum_{i=1}^{n_w} |w_i| = \frac{\lambda}{2m} \|w\|_1, \quad \|w\|_1 = \sum_i |w_i|,$$

### Dataset Augmentation

- Create fake data and add it the training set

### Noise Robustness

- the addition of noise with infinitesimal variance at the input of the model is equivalent to imposing a penalty on the norm of the weights
- For RNN
  - o noise has been used in the service of regularizing models by adding it to the weights
  - o a stochastic implementation of Bayesian inference over the weights
  - o Adding noise to the weights is a practical, stochastic way to reflect this uncertainty
- Injecting noise at the output targets
- assume that for some small constant  $\alpha$ , the training set label  $y$  is correct with probability  $1 - \alpha$ , and otherwise any of the other possible labels might be correct.
  - o label smoothing

### Dropout Regularization

- go through each of the layers of the network and set some probability of eliminating a node in NN
- Inverted dropout
  - o Randomly zero out different hidden units

### Multi-Task Learning

- improve generalization by pooling the examples (which can be seen as soft constraints imposed on the parameters) arising out of several tasks.
- when part of a model is shared across tasks, that part of the model is more constrained towards good values (assuming the sharing is justified), often yielding better generalization.

### Early Stopping

- We can obtain a model with better validation set error (and thus, hopefully better test set error) by returning to the parameter setting at the point in time with the lowest validation set error.
- Every time the error on the validation set improves, we store a copy of the model parameters.

- When the training algorithm terminates, we return these parameters, rather than the latest parameters.

## Optimization

### Normalization

- Normalize training sets

- Subtract mean

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

$$x := x - \mu$$

- Normalize variance

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

$$x := x / \hat{\sigma}$$

- Use same  $\mu$   $\hat{\sigma}^2$  to normalize test set

### Weight initialization for DN

- Partially solve Vanishing/exploding gradients problem
  - When training a very deep NN, the slopes can get either very big or very small, which makes training difficult
- Random initialization ([Hyperparameter](#))
  - Set the variance of  $W_i$  to be equal to  $1/n$

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

- ReLu: use  $2/n$
  - Tanh: use  $1/n$

### Gradient Checking

- kind of debugging your back prop algorithm
- Take  $W$  and  $b \rightarrow$  concatenate and reshape into a big vector  $\theta$

$$\mathcal{J}(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \mathcal{J}(\theta)$$

### Gradient Clipping

- Threshold the values of the gradients before performing a gradient descent step

### Minibatch

- Split examples to batches and each batch uses more than one but less than all of the training examples
- Mini-batch gradient descent
  - Process single batch of training samples at a time
- Take many gradient descent passes/epoch

### Bias correction in exponentially weighted averages

- make more accurate estimate during initial phase of the estimate

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

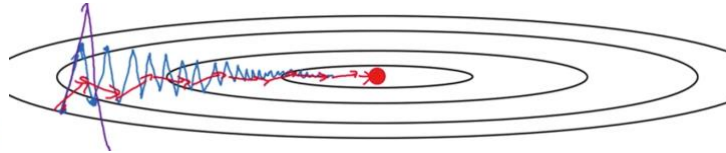
### Gradient descent with momentum

- compute an exponentially weighted average of your gradients and then use that gradient to update your weight
- take steps that are much smaller oscillations in the vertical direction and moving quickly to horizontal direction

$$V_{dw} = \beta V_{dw} + (1-\beta) \frac{dW}{dt}$$

$$V_{db} = \beta V_{db} + (1-\beta) \frac{db}{dt}$$

$$W := W - \alpha V_{dw}, \quad b := b - \alpha V_{db}$$



- Hyperparameters
  - o Learning rate  $\alpha$
  - o Control EWA:  $\beta$  (most common = 0.9)

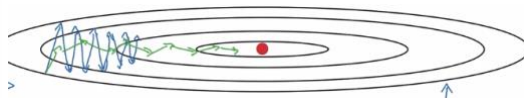
### RMSprop

- Speed up gradient descent
- Slow down the learning in the vertical direction and speed up learning in the horizontal direction

$$S_{dw} = \beta S_{dw} + (1-\beta) \frac{dW^2}{dt^2}$$

$$S_{db} = \beta S_{db} + (1-\beta) \frac{db^2}{dt^2}$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dw}}}, \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$



### Adam

- Idea: take Momentum and RMSprop, and put them together

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) \frac{dW}{dt}, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) \frac{db}{dt} \leftarrow \text{"momentum"} \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) \frac{dW^2}{dt^2}, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) \frac{db^2}{dt^2} \leftarrow \text{"RMSprop"} \beta_2$$

$$V_{dw}^{corrected} = V_{dw} / (1-\beta_1^t), \quad V_{db}^{corrected} = V_{db} / (1-\beta_1^t)$$

$$S_{dw}^{corrected} = S_{dw} / (1-\beta_2^t), \quad S_{db}^{corrected} = S_{db} / (1-\beta_2^t)$$

$$W := W - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}}, \quad b := b - \alpha \frac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$$

- Hyperparameters
  - o Learning Rate  $\alpha$
  - o  $\beta_1: 0.9$
  - o  $\beta_2: 0.999$
  - o  $\epsilon: 10^{-8}$

### Learning Rate decay

- Slowly reduce learning rate overtime
- Learning rate

$$\alpha = \frac{\alpha_0}{1 + \text{decay-rate} * \text{epoch-number}}$$

- o =
- o Decay-rate (hyperparameter)

$$\alpha = \alpha = 0.95^{\text{epoch num}} \cdot \alpha_0 \quad - \text{exponentially decaying}$$

## Hyperparameter Tuning

- choosing a set of optimal hyperparameters for a learning algorithm

Order	Hyperparameters
Most important	Learning rate
Less important	Momentum term ( $\beta = 0.9$ ); Mini-batch Size: # hidden units
Least Important	# of layers; Learning rate decay
Usually don't change	$\beta_1$ (0.9); $\beta_2$ (0.999); $\epsilon$ ( $10^{-8}$ )

## How to select a set of values to explore

- Coarse to fine finding scheme
  - o Zoom in a smaller region of the hyperparameters which work better and then sample more density within this space
- Appropriate scale for hyperparameters
  - o Search in a log scale
  - o Special case (Hyperparameters for EWA)

## Training many models in parallel

- Use different models with different set of the hyperparameters at the same time

## Batch normalization

- normalize the input layer by re-centering and re-scaling
- reduces the problem of coordinating updates across many layers
- Normalize any hidden layer (a) so as to train (W,b) faster
- Normalize (Z)

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta \quad \text{learnable parameters of model.}$$

$$x \xrightarrow{W^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\text{ReLU}(\text{Norm}(b))} \tilde{z}^{(1)} \xrightarrow{a^{(1)} = g(\tilde{z}^{(1)})} z^{(2)} \xrightarrow{W^{(2)}, b^{(2)}} z^{(2)} \xrightarrow{\text{BN}} \tilde{z}^{(2)} \xrightarrow{a^{(2)}} \dots$$

- Parameters: W, b,  $\gamma$ ,  $\beta$
- Working mini-batch

$$X^{(1)} \xrightarrow{W^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\text{BN}} \tilde{z}^{(1)} \xrightarrow{g^{(1)}(\tilde{z}^{(1)}) = a^{(1)}} z^{(2)} \xrightarrow{W^{(2)}, b^{(2)}} z^{(2)} \xrightarrow{\text{BN}} \tilde{z}^{(2)} \xrightarrow{\dots} \dots$$

$$X^{(2)} \xrightarrow{\dots} z^{(2)} \xrightarrow{\text{BN}} \tilde{z}^{(2)} \xrightarrow{\dots} \dots$$

- At test time

$$\mu = \frac{1}{m} \sum_i z^{(i)} \quad \sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \quad z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \leftarrow \mu, \sigma^2: \text{estimate using exponentially weighted average (across mini-batches)}$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

## Multi-class classification

### Softmax

- Classify more than 2 classes
- assigns decimal probabilities to each class in a multi-class problem. Those decimal probabilities must add up to 1.0

### Steps to build Deep learning model

1. set up dev/test metrics
2. Build initial ML System quickly
3. Use Bias, variance to prioritize next steps

## RNN

### Gated Recurrent Units

- A modification to the RNN hidden layer that makes it much better capturing long range connection

### LSTM

Peephole:

- Peephole connections allow the gates to access the constant error carousel (CEC), whose activation is the cell state.

CTC score function

- Used to find an RNN weight matrix that maximizes the probability of the label sequences in a training set, given the corresponding input sequences.