

# Gradient Descent and Optimization

Monday, December 13, 2021 9:16 PM

## 18. Gradient Descent and Optimization

### Optimization

- define cost func (objective func) to measure quality of soln
  - o minimize diff bw model prediction and labels
- use optimizer to search/solve for a soln that minimizes/maximizes cost func
  - o iteratively search for good solution

### Training NN Classifier

- objective: minimize diff bw label  $y$  and predicted label  $\hat{y}$
- prediction = func of input  $x$ , params  $w, b$
- training objective = func of params  $w, b$ , sample  $x$ , label  $y$
- for training set,  $x, y$  are const. (can exclude  $b$  for now)

$$J = f(y, \hat{y})$$

$$\hat{y} = h(x, w, b)$$

$$J = f(y, x, w, b)$$

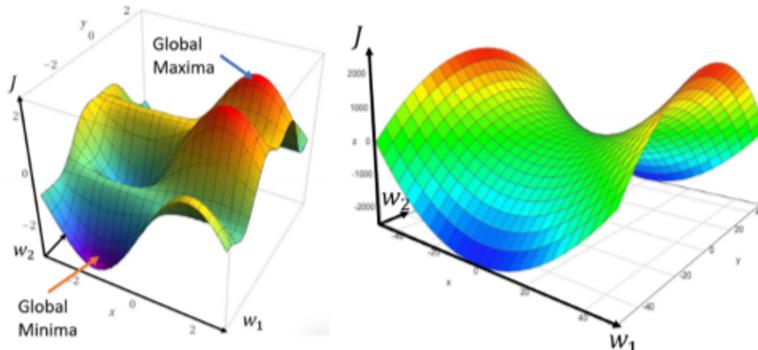
$$J = f(w, b)$$

$$J = f(w)$$

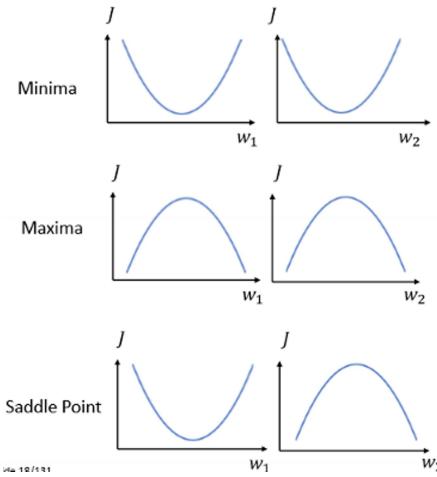
- for NN and training set, cost func  $J(w)$  is a func of only  $w$

### Global Optima

- o Optima: maxima and minima
- o Biggest/smallest among maxima/minima
- o local maxima/minima refers to the rest

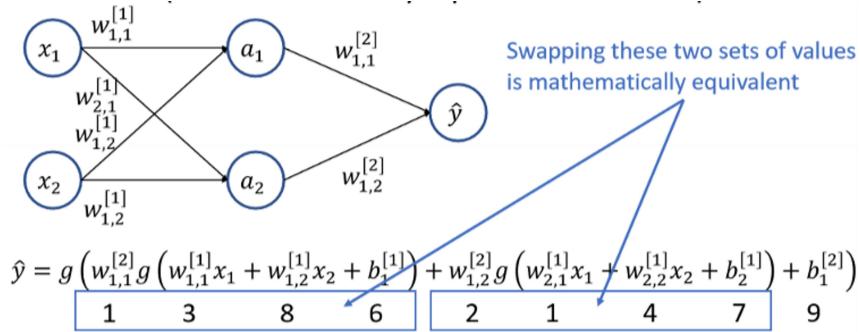


- gradients at optima and saddle points = 0,  $\frac{\partial J}{\partial w_i} = 0$  for all params  $w_i$



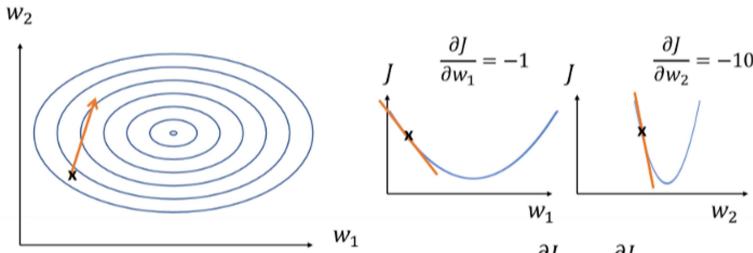
- DL Cost function is not convex

  - o many equal global minimal



## Gradient Descent

- start somewhere in space
- move in direction with steepest decrease in cost
- repeat
- Steps:
  - o Train model
  - o evaluate
  - o change hyperparam
  - o repeat
- **Problem:**
  - o too long to compute gradient for one training iteration
  - o requires too much memory to store activations
- **Mini-Batch Gradient Descent**
  - o use small subset of training set as approx of overall training set
  - o Sampling:
    - random shuffle full set
    - partition into mini-batch
    - iterate across each mini-batch
  - o one full pass through the set = epoch
  - o mini-batch size of 1 = stochastic gradient descent (SGD)
- **Problem:**
  - o **Different dimension change at diff rates**
    - direction of steepest descent isn't directly to min unless is a circle



Larger steps at steeper areas, and smaller steps at shallower areas

$$\frac{\partial J}{\partial w_1} \ll \frac{\partial J}{\partial w_2}$$

$$w_1 = w_1 - \alpha(-1)$$

$$w_2 = w_2 - \alpha(-10)$$

- Local optima and Saddle Points

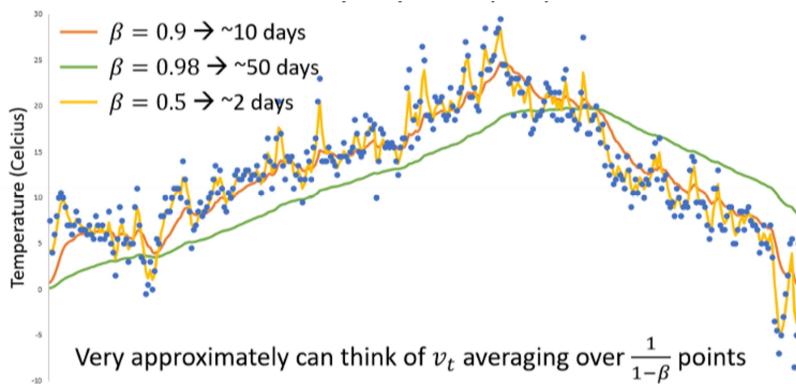
- gradient = 0 at local optima and saddle points
- saddle points are unstable but no gradient info
- gradient descent will stop updating params

- Meandering mini-batch gradients

### Gradient Descent with Momentum

- Use historical data with exp. weighted avg method to get out of 0 gradients saddle point
- Exponentially weighted averages

- moving average, smooth out fluctuations and highlight trends



- Solution:

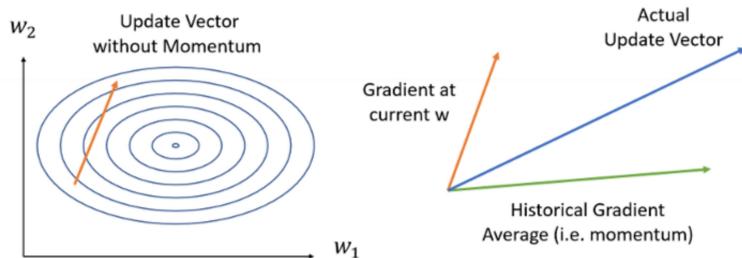
- instead of using gradient directly to update params, use exp weighted avg of past gradients

- Problem:

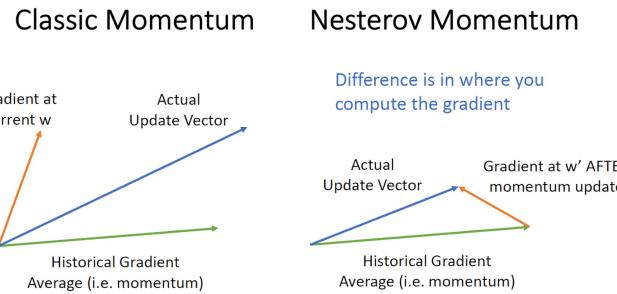
- Gradient in diff dimensions
  - consistent gradients build up velocity from accumulated acceleration
  - inconsistent gradients will cancel out
- Local minima and saddle point
  - at saddle point, gradient = 0, but momentum won't be (historical component)
  - at local minima, velocity can help get back out of local minima
- Meandering mini-batch gradients
  - moving avg will create smoothing effect

### Momentum so far

$$v_t = \beta \cdot v_{t-1} + \frac{\partial J}{\partial w}$$



## - Nesterov Momentum

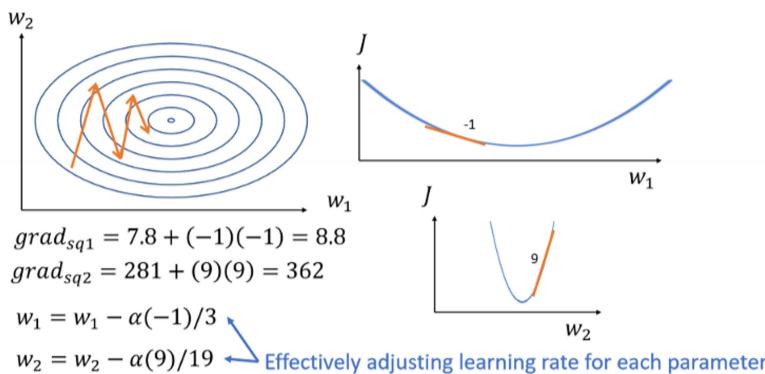


## Per-Parameter Adaptive Learning Rates

- to fix adaptive per parameter learning rates, use per-parameter learning rates

### - Adagrad

- o square of gradients focuses on magnitude, not direction
- o dampens steep gradients and accelerates shallow gradients by dividing by grad\_sq



### - Adam

- o combines RMSProp and Momentum
- o good default choice for optimizer

## Bias Correction for Exponentially weighted avg

- use biased vals in moving avg
- use corrected versions when doing gradient descent parameter update

$$v_{(t)}_{biased} = \beta \cdot v_{(t-1)}_{biased} + (1 - \beta) \cdot T_t$$

$$v_t = \frac{v_{(t)}_{biased}}{1 - \beta^t}$$

Denominator approaches 1 as t increases

- bias correction
  - o sorts itself out after several iterations but
  - o will make large updates at start which destroy weight initialization
  - o may result in no gradients

## Keras Optimizers:

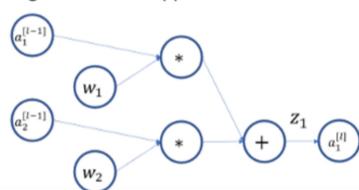
- SGD
- RMSProp
- Adagrad
- Adadelta
- Adam
- AdaMax
- Nadam

## Learning Rate Schedules

- how to choose learning rate: vary over time
    - o start high, reduce over time
    - o annealing/decaying learning rate
  - Decay learning rate
    - o reduce learning rate once progress plateaus
    - o Decay too slow = wasting time bounding
    - o Decay too fast = slow down training unnecessarily
  - decay based on func
    - Exponential Decay  $\alpha_t = \alpha_0 0.95^t$   $\alpha_t = \alpha_0 e^{-kt}$
    - Linear Decay  $\alpha_t = \alpha_0 \left(1 - \frac{t}{T}\right)$
    - Cosine Decay  $\alpha_t = \alpha_0 0.5 \left(1 + \cos\left(\frac{\pi t}{T}\right)\right)$
    - Inverse sqrt Decay  $\alpha_t = \alpha_0 \frac{1}{\sqrt{t}}$
    - 1/t Decay  $\alpha_t = \alpha_0 \frac{1}{1 + k * t}$
- $t$  - current training iteration
  - $T$  - Total training iterations
  - $\alpha_0$  - initial learning rate
  - $\alpha_t$  - learning rate at iteration  $t$
  - $k$  - a hyperparameter

## Parameter Initialization

- want gradients to be well behaved
  - o don't want to start off somewhere where gradients are all 0, not clear where to move
  - o don't start close to global minima, initializing with 0s
  - o initialize with a constant instead, **gaussian random dist** (breaks symmetry, not all initialized to same value)
  - o Gaussian with mean 0, can also expect that final weights might be zero-centered
  - o Gaussian random var \* x will give rand var a std dev = x
- **Gaussian random dist**

  - As you go deeper, activations of the units become 0 or close to 0
  - This means the gradients also approach 0
  - o for DNN, activation and gradients become 0
  - o can change std.dev  
*Was 0.01 before*  
 $W = 0.1 * np.random.randn(fan\_in, fan\_out)$
- need to find ideal variance
  - o Xavier initialization (for TanH)
  - o set variance = # inputs to layer (fan\_in)
 
$$W = (1 / np.sqrt(fan\_in)) * np.random.randn(fan\_in, fan\_out)$$
  - o Kaiming/He\_normal for ReLU

## Initializing Bias

- can initialize with 0s

## Gaussian vs Uniform Distribution

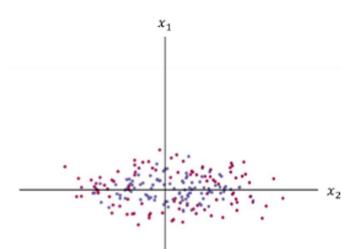
- some init schemes may draw from uniform dist instead of Gaussian

## Data Preprocessing

- sigmoid always positive
  - o all param updates will be positive
  - o inefficient training
  - o pick zero-centered activation func
  - o 1st layer, inputs are training data, no activation function -> non-zero centered?

- **Mean subtraction**

### Mean Subtraction



- Compute mean of each feature across all training samples

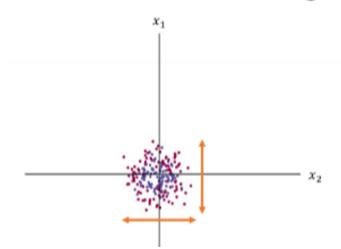
$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)}$$

- Subtract mean from each sample's features

$$x' = x - \mu$$

- **Normalization/Scaling**

### Normalization/Scaling



- Compute variance of each feature across all training samples

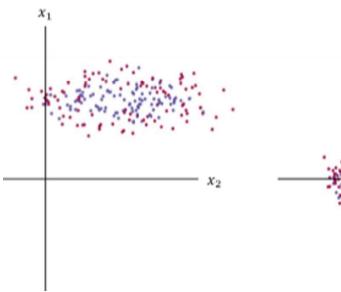
$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2$$

- Divide each feature by its standard deviation

$$x' = x / \sigma$$

- o relative feature scales, corresponding weights become same scale
- o absolute feature scales, don't want scales to be large, leads to large gradients, small change = big change in final cost, harder to optimize

- **Standardization (Z-Score Normalization)**



$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)}$$

$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2$$

$$x_i' = \frac{x_i - \mu_i}{\sigma_i}$$

## Image data

- each pixel = feature, same scale relative, but may still want to normalize
  - AlexNet
    - Subtracted mean image
    - i.e. computed mean on each feature (pixel) across dataset
  - VGGNet
    - Calculated a mean for each color channel (e.g.  $\mu_r, \mu_g, \mu_b$ )
    - Each channel's feature (pixel) subtracted its corresponding channel mean
  - ResNet
    - Subtracted channel mean like VGGNet
    - Also divided by channel standard deviation
- **At inference/prediction time**
  - o any transformation performed on input for training performed for inputs supplied at prediction time
  - o if used Z-Score Normalization, need to apply same func w/ same mean and std during

prediction time

$$x_i' = \frac{x_i - \mu_i}{\sigma_i}$$

### Batch Normalization

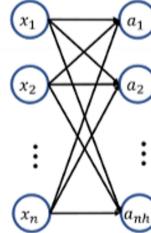
- Normalizing Network inputs
  - o helps training params in 1st layer, optimize
- normalizing inputs of hidden layers
  - o inputs to hidden layers are outputs of previous layers, batch normalization
  - o func of each layer changing as training progresses
  - o hard for intermediate layer optimization, inputs change
  - o small change = big change on deeper layer
  - o batch normalization helps stabilize optimization problem by giving layers target mean and variance

- For a given mini-batch with  $m$  samples,  $x$  is a matrix of shape  $(n, m)$

- For each input  $x_i$ , compute its

- mean  $\mu_i$  across the mini-batch  $\rightarrow n$  means
- variance  $\sigma_i^2$  across the mini-batch  $\rightarrow n$  variances

$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)} \quad \sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2$$



- For each sample and each feature, normalize  $\rightarrow (n, m)$

$$x_i' = \frac{x_i - \mu_i}{\sigma_i}$$

- forcing each layer's outputs to have zero-mean and unit variance too strict
- let model learn target mean and variance for each layer

### Learned mean and variance

- Two new trainable parameters,  $\gamma_i, \beta_i$  for each layer output that act to shift and scale the normalized layer outputs
- Normalize like before

$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)} \quad \sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2 \quad x_i' = \frac{x_i - \mu_i}{\sigma_i}$$

- Now shift and scale normalized values  $\rightarrow (n, m)$

$$\tilde{x}_i = \gamma_i x_i' + \beta_i$$

- if  $\gamma_i = \sigma_i$  and  $\beta_i = \mu_i$ , then  $\tilde{x}_i = x_i$  and hence  $\tilde{x}_i$  has the same mean and variance as the original input distribution
- if  $\gamma_i = 1$  and  $\beta_i = 0$ , then  $\tilde{x}_i = x_i'$  and hence  $\tilde{x}_i$  has 0 mean and unit variance
- For other values of  $\gamma_i$  and  $\beta_i$ ,  $x_i$  will have some other mean and variance

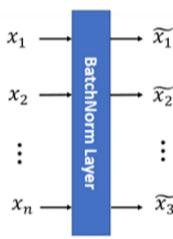
### Batch Norm layer

$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)}$$

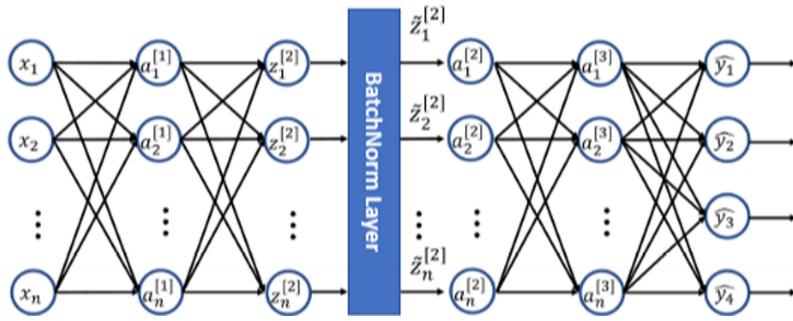
$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2$$

$$x_i' = \frac{x_i - \mu_i}{\sigma_i}$$

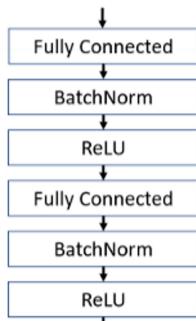
$$\tilde{x}_i = \gamma_i x_i' + \beta_i$$



- put layer before nonlinear activation

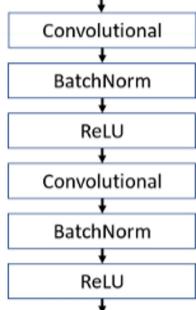


- BatchNorm and FCL



Term	Symbol	Shape
BatchNorm Inputs	$x$	$(n, m)$
Mean	$\mu$	$(n,)$
Stddev/Variance	$\sigma$	$(n,)$
Scale param	$\gamma$	$(n,)$
Shift param	$\beta$	$(n,)$
Normalized Inputs	$x'$	$(n, m)$
Batch Normalized Inputs	$\tilde{x}$	$(n, m)$

- BatchNorm and Conv Layer

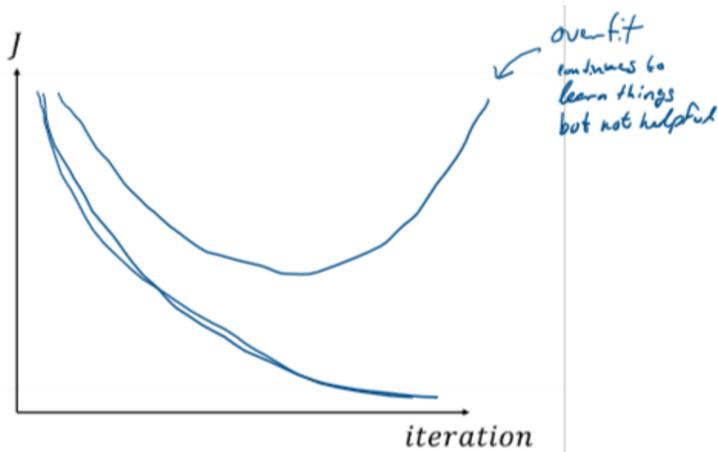


Term	Symbol	Shape
BatchNorm Inputs	$x$	$(h, w, c, m)$
Mean	$\mu$	$(c,)$
Stddev/Variance	$\sigma$	$(c,)$
Scale param	$\gamma$	$(c,)$
Shift param	$\beta$	$(c,)$
Normalized Inputs	$x'$	$(h, w, c, m)$
Batch Normalized Inputs	$\tilde{x}$	$(h, w, c, m)$

- At prediction time, keep moving average (exp weighted average), use these avg vals instead of  $\mu$  and  $\sigma$  during prediction time

## 20. Regularization and Training

### Overfitting



- model continues to learn things that does not help
- Need:
  - o more training data
  - o regularization

### Regularization via Cost func

- cost/obj func describes quality we want in our soln (param values)
- so far cost func only accounts for prediction loss
- can add terms to encourage regularization in our soln

$$J = \frac{1}{m} \sum_{j=1}^m L(\hat{y}, y) \quad J = \left( \frac{1}{m} \sum_{j=1}^m L(\hat{y}, y) \right) + R$$

Cost function so far

Cost function with  
Regularization Term

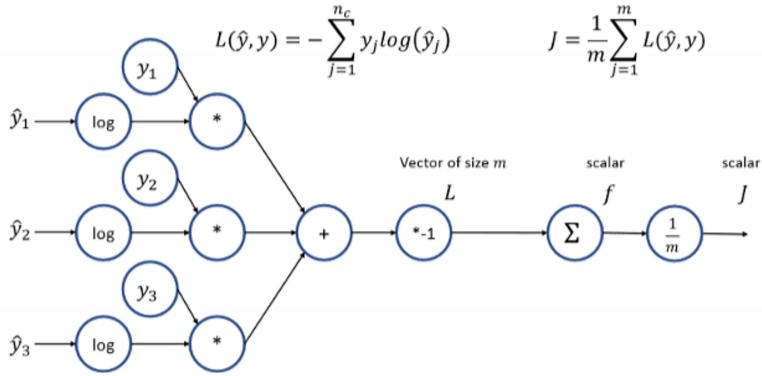
#### - L2 Regularization

- sum the square of each param value
- cost is minimized when each param value is small
- convex func independent of training set
  - regularization term is global min when all weights = 0

$$J = \frac{1}{m} \sum_{j=1}^m L(\hat{y}, y) + \sum w^2$$

- try to minimize the loss AND regularization terms
  - loss term = large if all weights = 0
- Since loss and regularization are competing terms, specify importance of each term based on lambda hyper param
  - lambda = 0: we don't optimize for regularization
  - lambda = inf: we don't optimize for loss
  - default in Keras = 0.01

$$J = \frac{1}{m} \sum_{j=1}^m L(\hat{y}, y) + \lambda \sum w^2$$



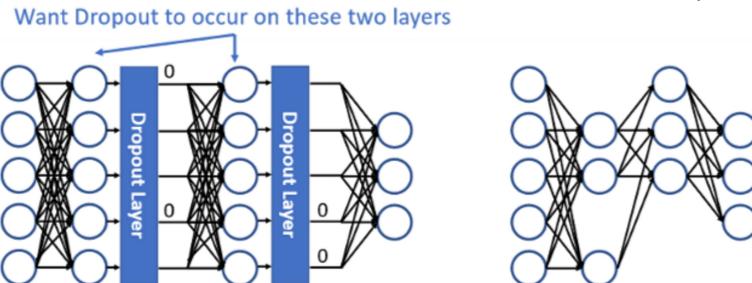
- Reduces overfitting by discouraging subset of weights dominating

### - Regularizing Bias Params

- doesn't have big impact since orders of magnitude more in weight params

## Dropout

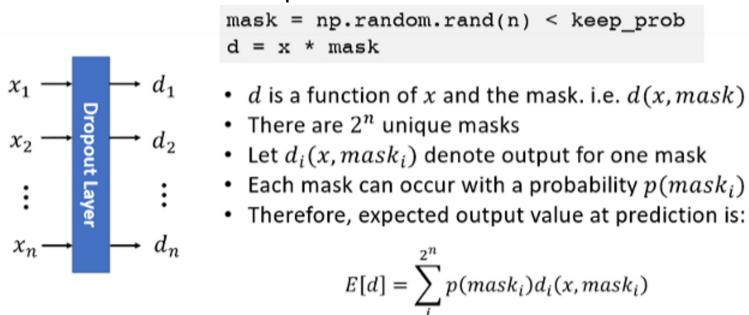
- on each param update iteration, **randomly** remove some hidden units from network
- training bunch of smaller simpler models and ensembling them together
- each model overfits in different ways, averages out
- don't put too much weight
- force each unit to learn to work well with rand subset of input units



Can implement Dropout by outputting 0  
at appropriate locations

### - At prediction time

- nondeterministic predictions!



- where to put dropout layers?
  - FCL since prone to overfitting compared to conv layers
  - conv layer less prone to overfitting since each swatch is a separate piece of training data

## DropConnect

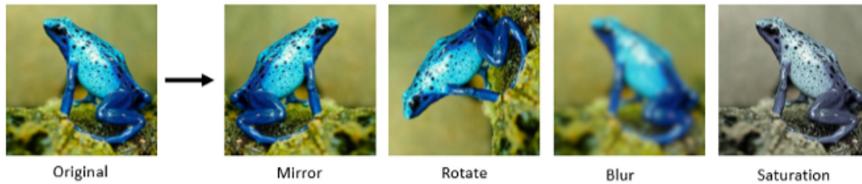
- Similar to dropout except 0 out random weights (connections) at training instead of nodes

## BatchNorm Regularization Effect

- mean and variance on minibatch is an approx to actual mean and variance
- introduces randomness
- unintended regularization effect

## Data Augmentation

- generating new training data from existing training data



- take different crops to generate new images

- **Types of regularization to use**

- L2 most common
- Dropout for large FCL
- don't rely on BatchNorm for regularization
- Data Augmentation for imaged data

## Hyperparameter Tuning

- Architecture
  - # layers
  - # units/filters per layer
- Optimization
  - learning rate
  - weight init
  - optimizer hyperparameters (momentum beta)
  - regularization techniques
- for some hyperparams, want to search over log scale, eg. learning rate range 0.0001 to 1
- coarse to fine
  - search in initial range of hyperparam values
  - find values that minimize cost
- **Training Advice**
  - start by using small subset of training set, get model to 100% accuracy (get training accuracy high first)
    - turn off regularization
    - flush out bugs in optimization flow
  - using full training set, find learning rate that shows good decrease in cost (close the gap and improve validation accuracy)
    - turn on regularization
    - see effect of learning rate in small number of training iters
  - do the hyperparam search as shown before
  - look at failing cases and cost curves (visualize)

## Looking at Cost Curves

\*\*\*Need to see lectures for curves graphs

## Transfer Learning

- take model trained for a task, repurpose it for second similar task
- keep some of the learnings (param values)
- used for img data and text/speech
- train new data set, but only let new output layer's params be updated
  - early layers of CNNs learn vocab of visual constructs
  - don't need to relearn these
- can let last couple conv layers also be retrained
- when to use
  - both tasks have same inputs
  - significantly less training data available for new task

- low level features are similar for tasks
- Pros
  - leverage previous training efforts, don't start from scratch
  - lets you start with good param values
  - don't need to relearn common low level features
  - helps pre train model if not much data