

## Linear Regression

Cost Function  $J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^i) - y^i)^2$  prediction model

Gradient Descent repeat:  $w = w - \alpha \frac{\partial}{\partial w} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^i) - y^i) x_j^i$

"bias" / intercept value  $b = b - \alpha \frac{\partial}{\partial b} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^i) - y^i)$   $j$ th feature

Gradient Descent();

Cost Function (Model Equation)

$x^{(i)}$   $i$ th training entry

$\vec{x}$  now vector of feature values

$x_j$  =  $j$ th feature (column vector)

$x_j^{(i)}$  =  $j$ th feature in vector of  $i$ th training entry

$$\begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix} \cdot \begin{bmatrix} \cdot & \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \end{bmatrix}$$

$$f_{w,b}(\vec{x}) = \vec{x} \cdot \vec{w} + b = [w_1, w_2, w_3, \dots, w_n] \cdot [x_1, x_2, x_3, \dots, x_n] + b \in w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + b$$

## Vectorized Multiple Linear Regression

Gradient Descent:

Loop, adjust weights:

$w = w - \underbrace{\# \text{Compute gradient}()}_{\substack{\text{loop each } X \text{ training entry;} \\ \text{updates} \\ \text{weights \& bias}}}$  (vector product  $f(x)$  prediction -  $y$  training entry).

$\underbrace{\text{loop each } j \text{ feature:}}_{\substack{\text{loop each } j \text{th } X \text{-entry feature} \\ \text{add } * \text{ jth } X \text{-entry feature}}} \rightarrow \text{gets vector of } \underbrace{\begin{bmatrix} \frac{\partial}{\partial w_1}, \dots, \frac{\partial}{\partial w_n} \end{bmatrix}}_{\substack{\text{of slopes} \\ \text{size } n}}$

$$b = b - \dots$$

Feature Scaling - used to normalize the scaling of parameters to optimize gradient descent convergence always types: Mean Normalization:  $x_i - \bar{x}_i$  Z-Score:  $\frac{x_i - \bar{x}_i}{\sigma_i}$   $\bar{x}_i$  mean (mu)  $\sigma_i$  standard deviation (scale factor)

$$\text{sum for } -1 \leq x_i \leq 1$$

normalize inputs after creating  $w$  weights.

normalizing doesn't change cost b/c its proportional (both  $X$  and  $w$  are scaled)  $\sigma_i$  standard deviation (variance)  $\sigma_i$  sigma

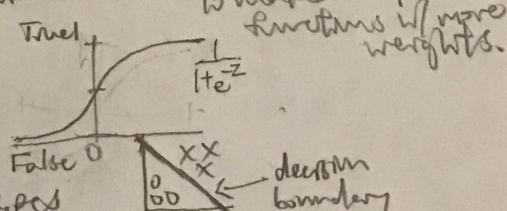
normalize  $X$  and  $w$ , NOT  $y$  training data - will be converted to correct scale.

Feature Engineering - combining features to make new ones

Polynomial Regression:  $[x, x^2, x^3, \dots, x^n]$  one parameter can be used for others degrees, to make complex functions w/ more weights.

## Logistic Regression (Classification)

Sigmoid function:  $g(z) = \frac{1}{1+e^{-z}}$ ,  $z = \vec{w} \cdot \vec{x} + b \rightarrow = \frac{1}{1+e^{(w \cdot x+b)}}$



Decision Boundary: Line of model between different classification types

## Logistic Cost Function

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\vec{w}, b}(\vec{x}^i), y^i) \rightarrow \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^i)), & \text{if } y^i = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^i)), & \text{if } y^i = 0 \end{cases}$$

$$- J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m (y^i \log(f_{\vec{w}, b}(\vec{x}^i)) + (1-y^i) \log(1 - f_{\vec{w}, b}(\vec{x}^i))) \rightarrow = -y^i \log(f_{\vec{w}, b}(\vec{x}^i)) - (1-y^i) \log(1 - f_{\vec{w}, b}(\vec{x}^i))$$

Gradient Descent();

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^i) - y^i) x_j^i, f_{\vec{w}, b}(\vec{x}) = \frac{1}{1+e^{(\vec{w} \cdot \vec{x} + b)}}$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^i) - y^i)$$

Cost function not used in gradient descent only for tracking progress

Regularization - for any model, you want it to generalize the accuracy of new input  $\vec{x}$ .  
 changes values of  $w$  Problem B: if it's too specific for data, it overfits. If too vague, it underfits.

Regularization creates a balance

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \left[ (f_{\vec{w}, b}(\vec{x}^i) - y^i)^2 + \frac{\lambda}{m} \sum_{j=1}^n w_j^2 \right]$$

minimize  $\xrightarrow{\text{mean squared error}}$  the term balances  $\xrightarrow{\text{regularization}}$  the tradeoff.

Gradient Descent, Regularized Linear Regression, same for Logistic

Repeat:

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b) \rightarrow \left[ \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^i) - y^i) x_j^i + \frac{\lambda}{m} w_j \right]$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^i) - y^i)$$

$$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2^2 + w_3 x_3^3 + b$$

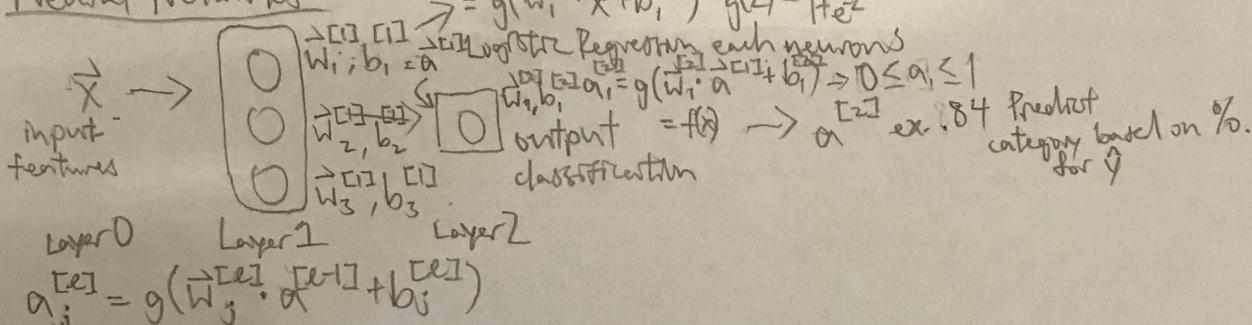
reduce  $w_j$  to reduce overfitting

## Regularized Logistic Regression

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[ y^i \log(f_{\vec{w}, b}(\vec{x}^i)) + (1-y^i) \log(1-f_{\vec{w}, b}(\vec{x}^i)) \right] + \frac{\lambda}{m} \sum_{j=1}^n w_j^2$$

To address overfitting, try to ① Add more training points, ② Apply regularization, ③ Reduce # of features

## Neural Networks



Forward Propagation - computations "forward" left to right

Sequential() : create neural net w/ layers and weights  $W$  and  $b$

Dense() : create a layer

$$A = g(A^T \cdot W + B)$$

4 params

$$\begin{matrix} A^T \\ \vec{w}_1, b_1 \\ \vec{w}_2, b_2 \\ \vec{w}_3, b_3 \end{matrix} \quad \begin{matrix} [1, 2], 3 \\ [4, 5], 6 \\ [7, 8], 9 \end{matrix} \quad \left\{ \begin{matrix} A^T \cdot W + B \text{ sigmoid} \\ A^T = [1, 2] \cdot \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \end{bmatrix} + [3, 6, 9] \text{ sigmoid} \end{matrix} \right.$$

$$1 \times 2 \cdot 2 \times 3 + 1 \times 3 = 1 \times 3$$

Activation Function Types: Binary Classification

$$y = \frac{1}{1+e^{-z}}$$

Sigmoid

Regression

$$y = z$$

Linear

$$y = \frac{1}{1+e^{-z}}$$

Regression

$$y = \frac{1}{1+e^{-z}}$$

ReLU

$$y = \max(0, z)$$

rectified linear unit

Use for hidden layers

## Softmax Regression - multiclass classification

for  $j, 1-N$  outputs,

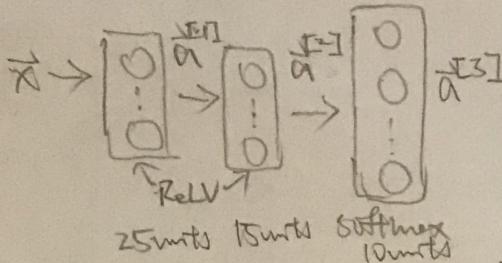
$$z_j = w_j \cdot \vec{x} + b_j$$

$$a_j = \frac{e^{z_j}}{\sum_k e^{z_k}}$$

$$\text{loss}(a_1, \dots, a_N, y) = \begin{cases} -\log a_j & \text{if } y=j \\ -\log a_j & \text{if } y \neq j \\ -\log a_N & \text{if } y=N \end{cases}$$

$$\rightarrow P(y=j|\vec{x}) = \frac{e^{z_j}}{\sum_k e^{z_k}}$$

## Softmax Output



$$z_i = w_i \cdot x + b_i$$

$$a_i = \frac{e^{z_i}}{e^{z_1} + e^{z_2} + e^{z_3}} = P(y=i|x)$$

$$z_{10} = w_{10} \cdot x + b_{10}$$

$$a_{10} = \frac{e^{z_{10}}}{e^{z_1} + \dots + e^{z_{10}}} = P(y=10|x)$$

Need to get  $\sum_k z_k$  first to use for each  $a_k$

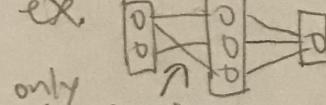
More numerically accurate version:  $L(\vec{a}, y) = \begin{cases} -\log \frac{e^{z_y}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y=1 \\ \text{softmax loss value} & \text{if } y \neq 1 \end{cases}$  Put on one line, put softmax into loss function

Multilabel classification - not the same as

Optimization multilabs. Inputs can have multiple outputs if  $y=10$  function  $\log z$ , will be more accurate of gradient descent where each parameter  $w_i$ ,  $b_i$  has a different  $\alpha$  (learning rate)

Dense layer:

Convolutional layer: activation units takes in only select activation units from previous layer - leads to faster computation and reduces overfitting.



Back Propagation: the use of derivatives to find the weights of the neurons in a neural network, first feeding forward, then getting the partial derivative of each term in the cost function in a computation graph with respect to the intermediate value, propagating left to right to find out how much to nudge the weights for each neuron and

Train not all examples, keep ~10-20% to test your model, after training model, test the cost by (how well it is) getting cost of unused test examples and compare with Get cost, split data into 3 parts of cross-validation set, then choose best training cost value cost and use untouched test set to test the generalization error?

## ML Advice

Train ~60-70% of data, cross-validation set test which model to use, then use remaining test set to determine how good the model is (3 sets).

1. Split data into training set, cross-validation set, and test set
2. Try many different NN architectures (hyperparameters) and train on train set (scale it)
3. Use model and find cost of cross-validation set
4. The model with the lowest cv cost is the most appropriate model that keeps the cost low while generalizing to new examples.
5. Test model on test set to get accuracy.

## Bias and Variance (Over & Underfitting)

High bias (underfit)

Jtrain high

Jcv high

High variance (overfit)

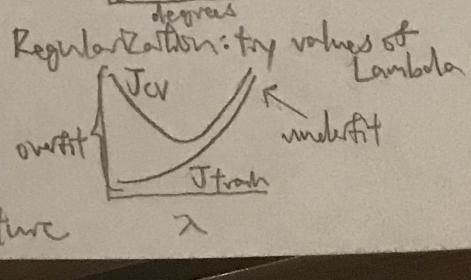
Jtrain low

Jcv high

Just right:

Jtrain low

Jcv low



To better the model, try getting data, changing # of features, changing  $\lambda$ , changing model architecture

Compare  $J_{\text{train}}$  and  $J_{\text{cv}}$  to baseline level, that is, how well you expect the algorithm to perform (would be how well humans can do task or other algos). Based on the difference of  $J_{\text{train}}$  and  $J_{\text{cv}}$  to the baseline, then determine if it is overfitting (high variance) or underfitting.

Plot graphs of different models vs. costs and lambda values vs. cost to see if it has high bias or variance. Increasing the size of neural networks increases performance, need to regularize, it is more computationally expensive, NN's usually have high variance. Error Analysis - look at examples the model misclassified to diagnose problems and certain data that can help from error analysis.

Data Augmentation - changing/distorting a training example slightly to make more examples and improve the examples and improve the algo.

Data Synthesis - creating artificial data

Transfer Learning - use NN parameters from another project (combined) and use those initial weights for model - teaches the model to detect similar patterns

Precision and Recall - if dataset is skewed (more examples that output a certain value), precision/recall table is useful relative to when our model was correct vs when we predicted that class and was wrong, and correctness vs. when we predicted another class but it was

Precision: model predicted a class, how accurate were these predictions? Actually, the <sup>(1)</sup> wrong

Recall: ground truth label was this class, how accurate were we for this label?

$\text{precision} = \frac{\text{true positives}}{\text{all predicted positives}}$        $\text{recall} = \frac{\text{true positives}}{\text{all actual positives}}$  Often, there is a tradeoff between precision and recall

optimal threshold = .99  
Recall = .01

"Harmonic Mean"  $\frac{2PR}{P+R}$        $P = \text{precision}$ ,  $R = \text{recall}$

To combine precision and recall, use F1 Score =  $\frac{1}{\frac{1}{2}(\frac{1}{P} + \frac{1}{R})} = \frac{2PR}{P+R}$  for the case of balancing the two highest

Regression: compare MSE cost of  $J_{\text{train}}$  and  $J_{\text{cv}}$

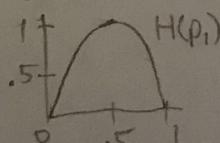
Classification: compare % accuracy of  $J_{\text{train}}$  and  $J_{\text{cv}}$

Decision Trees: a tool, tree-like structure that classifies predictions from a set of conditions

A decision tree must generate many decision nodes (how to split features) to split and how many. - splits to get pure batches (all 1 class) features to split.

Entropy: measure of impurity of data

$P_i = \text{fraction of examples with common class}$   $p_i = 1 - p_i$

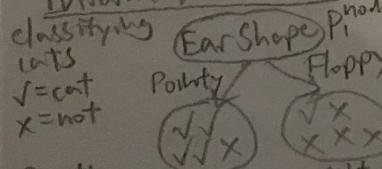


$$H(p_i) = -p_i \log_2(p_i) - (1-p_i) \log_2(1-p_i)$$

$$\begin{aligned} & \text{50-50 split } H(0.5) = 1 \quad \text{very pure} \quad H(0) = 0 = H(1) \\ & \text{has entropy} = 1, \quad \text{splits } 100-0 \\ & \text{or } 0-100, \text{ entropy} = 0 \end{aligned}$$

$$H\left(\frac{1}{2}\right) = 0.65 \quad \text{"not that impure"}$$

Information Gain: how much a split reduces entropy (purifies data) choose feature split with highest information gain



fraction in left subtree =  $4/5$        $p_i^{\text{left}} = 4/5$   
 fraction of node left =  $4/5$        $w^{\text{left}} = 5/10$   
 fraction of node right =  $1/5$        $w^{\text{right}} = 5/10$

$$\begin{aligned} \text{Information Gain} &= H(p_i) - (w^{\text{left}} H(p_i^{\text{left}}) + w^{\text{right}} H(p_i^{\text{right}})) \\ &\text{Weighted entropy} \end{aligned}$$

Ex One Hot Encoding - split up one feature with  $k$  values into  $k$  features with 1 value 0 or 1.

Ear Shape	Round	Pointy	Oval	
Round	1	0	0	Can also be used
Pointy	0	1	0	for Neural Nets
Oval	0	0	1	

Continuous valued features: choose a threshold to split data that maximizes information gain.

Regression Trees - decision trees with continuous y-labels.

Ex. Guessing weights of animals

Choosing a Split

Ear Shape		Variance root node: 20.51
Pointy	Floppy	
5	5	
Weights: 7.2, 9.2, 8.4, 7.6, 10.2	Weights: 8.8, 15, 11, 14, 6, 20	
Variance: 1.47	Variance: 21.87	
$W_{\text{left}} = \frac{5}{10}$	$W_{\text{right}} = \frac{5}{10}$	# in subtree # in root node

Regression Information

Gain

$$\text{Variance} - (w_{\text{left}}(\text{variance}_{\text{left}}) + w_{\text{right}}(\text{variance}_{\text{right}}))$$

$$20.51 - \left( \frac{5}{10}(1.47) + \frac{5}{10}(21.87) \right)$$

$$= 8.84$$

Stop splitting when chosen criteria: depth of tree, information gain threshold, # of examples in node, purity of node, or combination

Sampling with Replacement: make training set by choosing random values (keep them in set after choosing). To make multiple, varying training sets for ensemble of trees - duplicates are expected

Random Forest Algorithm:

For  $b=1$  to  $B$  decision trees:

Make "new" training set by sampling with replacement all  $m$  examples

Train decision tree on new dataset

To make prediction, run it through all generated decision trees and get average output

For Random Forest Algorithm, when choosing split feature, limit options to  $k < n$  (try  $k = \sqrt{n}$ )

Bagged Decision Tree - same thing but all features can be split randomly to diversify different decision trees

XGBoost Algorithm - even better decision tree algorithm that prioritizes making decision trees

NNs VS. Decision Trees

Neural Networks

Good for Structured & Unstructured data

Slow training

Block for misclassified examples

Box also built in regularization

A bit more interpretable and stop criteria, etc

Decision Trees

Only good on structured data

Faster training

Clustering - unsupervised learning technique for grouping data into categories

K-means Clustering Algorithm - gets "center of mass/cluster" of k cluster centroids

Randomly initialize k cluster centroids  $M_1, M_2, \dots, M_K$

Repeat until convergence:

Get distance of all points to cluster centroids, assign the cluster w/  
loop all points  $x_i$ :  $c^i = \text{cluster closest to } x^i$

cluster closest to point  $x^i$ ,  $= c^i$  ( $1 \text{ to } K$ ) which is  $\min_k (x^i - M_k)^2$

Move all clusters to center point of points closest to it. ( $1 \text{ to } K$ ) for each centroid 1 to  $K$ :

Cost Function (Distortion function)  $J(C, M) = \frac{1}{m} \sum_{i=1}^m \|x^i - M_{c^i}\|^2$   $\rightarrow$  sum of points in cluster over # of points in cluster

$M_k = \frac{1}{\|C_k\|} \sum_{i \in C_k} x^i$   $\rightarrow$  average in cluster  
 $C^i = \text{index of cluster } (K \text{ clusters}) x^i \text{ is assigned}$

K-means optimizes cluster distance  
centroids to minimize distortion cost

K-means Random Initialization - initialize cluster centroids at random examples and run H, getting lowest distortion cost

For  $i=1$  to  $100^{\leq 50-100}$ : ← Iterations of K-means

Randomly initialize K cluster centroids at the training examples

Run K-means to convergence, get  $C^1 \dots C^m$  (cluster  $i$ th example belongs to) and  $M_1 \dots M_K$  (the cluster centroid)

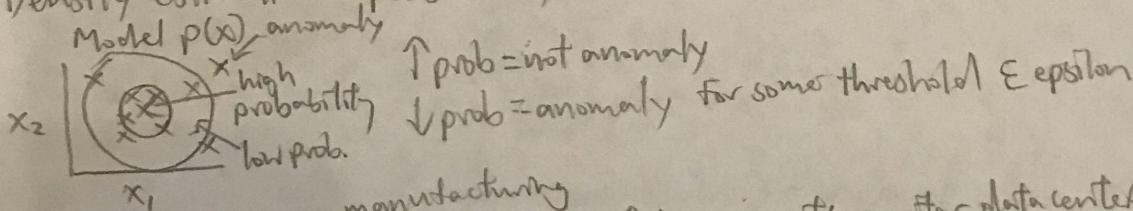
Compute cost function, store it  $\rightarrow J(C^1 \dots C^m, M_1 \dots M_K)$

\* Pick sets of clusters w/ lowest cost

The number of clusters  $K$  to use is ambiguous, can use elbow method, but usually depends on the purpose of the program

Anomaly Detection - algorithm to detect unusual inputs deviating from what is expected & tradeoffs.

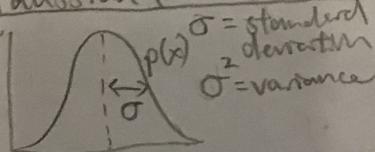
Density estimation - probability that input is from distribution

Model  $p(x)$  anomaly  
  
 $x_1$   $x_2$   $x$  ↑ prob = not anomaly  
↑ prob = anomaly for some threshold  $\epsilon$  epsilon

Examples:

Fraud detection, product inspection, fake accounts, monitor data centers

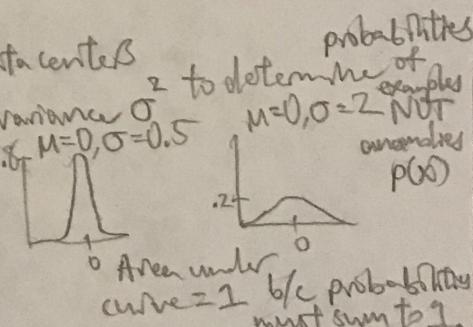
Gaussian (Normal) Distribution - bell shaped curve with mean  $M$ , variance  $\sigma^2$



$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-M)^2}{2\sigma^2}}$$

constants  $M = \text{mean}$   $\sigma = \text{std deviation}$

$M$  changes horizontally shift  
 $\sigma$  changes width



Area under curve = 1 b/c probabilities must sum to 1.

How to determine  $M$  and  $\sigma$ :

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x^i \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^i - \bar{x})^2 \quad \text{Called maximum likelihood estimation}$$

with  $n$  features

Density Estimation - get probability that an example is "normal"

Training set:  $m \times n$   $[x^1, x^2, \dots, x^m]$

$p(\vec{x}) = p(x_1; M_1, \sigma_1^2) \cdot p(x_2; M_2, \sigma_2^2) \cdots p(x_n; M_n, \sigma_n^2) \rightarrow$  probability of each feature across all examples in being "normal" multiplied together  $\rightarrow$  probability the examples are "normal"

$$\text{Distribution for example with } n \text{ features} = \prod_{j=1}^n p(x_j; M_j, \sigma_j^2)$$

## Anomaly Detection Algorithms

1. Choose  $n$  features for training example that might show anomalies
2. Fit parameters  $\vec{M}$  ( $1 \times n$ ) and  $\vec{\sigma}^2$  ( $1 \times n$ )

$$\vec{M} = \frac{1}{m} \sum_{i=1}^m \vec{x}^i \rightarrow [M_1, M_2, \dots, M_n] \quad \vec{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (\vec{x}^i - \vec{M})^2 \rightarrow [\sigma_1, \sigma_2, \dots, \sigma_n]$$

3. Compute  $p(\vec{x})$

$$p(\vec{x}) = \prod_{j=1}^n p(x_j; M_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - M_j)^2}{2\sigma_j^2}\right) \rightarrow \begin{matrix} \text{probability of 1 example not being} \\ \text{an anomaly} \end{matrix}$$

4. If  $p(\vec{x}) < \epsilon$ , is anomaly

Algorithm pushes heavily  $p(\vec{x})$  if any features have low probability b/c it multiplies all  $n$  prob.

Real-number evaluation - test how well the algorithm is working making  $p(x)$  ↓.

If you have some labeled data of anomalous and non-anomalous, do this:

heavy pushed  
split into 3 sets  
in CV set,

Training set:  $[\vec{x}^1, \dots, \vec{x}^m]$  assume normal examples

Cross validation set:  $[(\vec{x}_c^1, y_c^1), \dots, (\vec{x}_c^m, y_c^m)]$  } include same # of time  $\epsilon$  and  $X_j$  features

Test set:  $[(\vec{x}_{\text{test}}^1, y_{\text{test}}^1), \dots, (\vec{x}_{\text{test}}^m, y_{\text{test}}^m)]$  } some anomalies examples, labeled.

Alternative: no test set, just training set and CV set ← with anomalies, time  $\epsilon$  and  $X_j$ .

Use if we don't have many anomaly examples - can't test, may overfit

- Consider applying evaluation metrics

Anomaly Detection Vs. Supervised Learning

- Small # positive examples Large # positive and negative examples

- Algo detects examples similar to what it has seen before, good if next example anomalies are expected to be similar

Ex: Spam - similar

Common spam emails

manufacturing detects scratches

\* Predict weather (sunny, rainy, etc.)

\* Disease classification (known diseases)

Precision =  $\frac{\# \text{predicted 1 was right}}{\# \text{predicted 1}}$  b/c dataset is often skewed w/ anomaly detection

Recall =  $\frac{\# \text{predicted 1 was right}}{\# \text{ground truth label y was 1}}$

$$F1\text{-Score} = 2 \frac{P \cdot R}{P+R} \text{ prediction accuracy } [0,1]$$

## Choosing Which Features to Use

Distribution of values for a feature

It is helpful to make the features more Gaussian (bell-shaped) by using log or some function  $f(x_j)$

Error analysis - look at incorrect examples, change features (or add), or combine features

$$N_m \times N_u$$

## Recommender Systems

Movie	Abel(1)	Bob(2)	Carol(3)	Dave(4)	$x_1$ (romance)	$x_2$ (action)	$y_{(i,j)}$
Love at Last	5	5	0	0	0.9	0	
Romance Forever	5	?	?	0	1.0	0.01	
Cute puppies & love	?	4	0	?	0.99	0	
Nonstop car chases	0	0	5	4	0.1	1.0	
Swords vs. Karate	0	0	5	?	0	0.9	

Predict user j's rating for movie i as  $w_i \cdot x^i + b_i$

$$x^i = 1 \times n \quad w = n \times 1$$

$N_m = 5$  #movies  
 $N_u = 4$  #users  
 $n = 2$  #features  
 $r(i,j) = 1$  if user j has rated movie i (otherwise 0)  
 $y_{(i,j)} =$  user j rating on i  
 $w^i, b^i =$  user j parameters  
 $x^i = [0.99 \ 0]$   $w^i = [5 \ 0]$   $b^i = 0$   $x^i$  = feature vector for movie i  
 $w^i \cdot x^i + b^i = 4.95$   $m^i =$  #movies user j rated  
 Goal: learn  $w^i, b^i$

## Cost Function for user $j$ given features get weights

$$J(\vec{w}, b) = \frac{1}{2} \sum_{i:r(i,j)=1}^m (\vec{w} \cdot \vec{x}_i + b - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (w_k)^2$$

choose  $w$  and  $j$  to minimize cost b/w actual rating and predicted

Learn  $w$  for all users: only rated movies

$$J(\vec{w}^1, \dots, \vec{w}^m, b^1, \dots, b^n) = \frac{1}{2} \sum_{j=1}^n \sum_{i:r(i,j)=1} (\vec{w} \cdot \vec{x}_i + b - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^m (w_k)^2$$

If we don't have features  $x_1 - x_n$ , we can find them if we know weights  $w$  and  $b$ .  
ex.  $w_i \cdot x_i + b \rightarrow x' = [1 \ 0]$  find features for each movie

## Cost function for 1 movie $i$ given weights, get

$$J(\vec{x}_i) = \frac{1}{2} \sum_{j:r(i,j)=1}^m (\vec{w} \cdot \vec{x}_i + b - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (w_k)^2$$

$$J(\vec{x}_1, \dots, \vec{x}_m) = \frac{1}{2} \sum_{i:j:r(i,j)=1}^m (\vec{w} \cdot \vec{x}_i + b - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (w_k)^2$$

Collaborative Filtering - data of multiple users ratings to predict features for movies,  $w$  and  $b$ , must learn more features and user parameters to predict ratings

## Cost put together:

$$J(w, x, b) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} (\vec{w} \cdot \vec{x}_i + b - y^{(i,j)})^2 + \underbrace{\frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^m (w_k)^2}_{\text{regularization}} + \underbrace{\frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (w_k)^2}_{\text{regularization}}$$

## Gradient Descent

repeat:

$$\begin{aligned} w_i &= w_i - \alpha \frac{\partial}{\partial w_i} J(w, b, x) && \text{parameters} \\ b^i &= b^i - \alpha \frac{\partial}{\partial b^i} J(w, b, x) && \text{are } w, x, b \\ x_k^i &= x_k^i - \alpha \frac{\partial}{\partial x_k^i} J(w, b, x) \end{aligned}$$

Binary Labels - use sigmoid function to user liked it, binary classification

Ratings mean 1-user engaged, 0-not engaged,  $J$ -not shown

$$f_{(w,b)}(x) = g(w \cdot x + b), g = \frac{1}{1+e^{-z}}$$

$$\text{Binary Loss: } L(f_{(w,b)}(x), y^{(i,j)}) = -y^{(i,j)} \log(f_{(w,b)}(x)) - (1-y^{(i,j)}) \log(1-f_{(w,b)}(x))$$

$$J(w, b, x) = \sum_{(i,j):r(i,j)=1} L(f_{(w,b)}(x), y^{(i,j)})$$

Mean Normalization - needed b/c new user would have ratings 0. Want "natural" rating prediction

change  $y \leftarrow y - M \leftarrow \text{row mean}$  So, for new user with  $w=0=b$ ,

changes  $f_{(w,x,b)}(x) = w \cdot x + b + M$ ; with new  $f_{(w,b)}(x)$ , predictions =  $M \rightarrow$  average of other users' ratings for new user's natural ratings.

so prediction is same

\*compute column mean in some cases

Finding Related Items - how related are items to one?

Got feature vectors, smallest distance between item  $k$  and item  $x_i$

$$\sum_{k=1}^n (x_k^i - x_i^k)^2 \text{ also } \|x^k - x_i\|^2 \text{ smallest val is most related to } x_i$$

## Limitations

- tough for new items w/ few user ratings \* cold start - also for new user hard to recommend
- doesn't use additional site information about user (only the rating)

## Collaborative Filtering

- Recommend items to you based on other users' ratings

VS.

## Content-based Filtering

- Recommend items to you based on user features and item (including ratings) to find good matches

### Content-based Filtering

$X_u$ : user features for user  $j$  ex. age, movies watched, ratings, ...  
 $X_m$ : movie features for movie  $i$  ex. year, genre, reviews, ratings

Predict rating user  $j$  on movie  $i$ :

$$v_u \cdot v_m$$

shape of  $v_u$  and  $v_m$  must be same

$\uparrow$  computed from  $X_u$  features  
 computed from  $X_m$  features

= convert user & movie feature into vectors w/ same shape

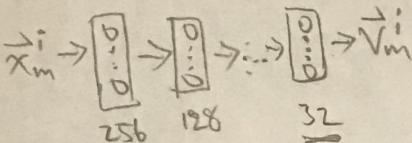
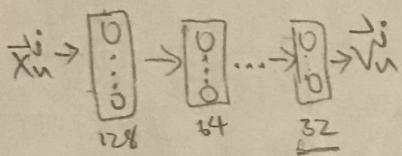
Getting  $v_u$  and  $v_m$  vectors  
 $X_u \rightarrow v_u$  "Web" network

- use neural networks "deep learning"

$X_m \rightarrow v_m$  "movies/items network"

User  $j$

Movie  $i$



$$\text{Prediction: } v_u \cdot v_m$$

can use sigmoid classification:  $g(v_u \cdot v_m)$  too

### Cost function

$$J = \sum_{(u, i) \in r(u)} (v_u \cdot v_i - y^{(u, i)})^2 + \text{NN regularization term}$$

$v_u$  is vector length 32, represents user  $j$  w/ features  $X_u$

$v_i$  is vector length 32, represents movie  $i$  w/ feature  $X_m$

Finding similar items:  $\|v_m - v_i\|^2$  small distances more similar

\* This can be pre-computed (overnight!), so each movie has similarity score for each movie ( $N_u \times N_m$ )

Recommendation moves from large catalogue: ① Retrieval, ② Ranking for quick access to the top

Retrieval - generate large list of possible items \* trade off w/ more items  $\Rightarrow$  performance vs. speed

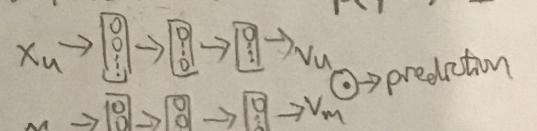
- ex prime out items 1. For last 10 movies user watched, pick 10 similar movies for each  $\Rightarrow \|v_m - v_i\|^2$   
 2. Pick 10 movies from user's 3 most viewed genres  $\Rightarrow$  analyze how many retrieved items to have pre-computed?

3. Top 20 movies in country

- then filter list; remove duplicates, items already watched

high prob. recommendations  $p(y^{(u, i)}) = 1$

Ranking - use the list (~100 items) and rank items with model



- order by highest rating

\* can pre-compute  $v_u$  and  $v_m$  for each user and movie and just do  $v_u \cdot v_m$  to get prediction quickly

### Ethics

Different ways to configure recommender systems:

- movies that user will rate higher
- products user will purchase
- ads optimizing for clicks
- generating most profit
- optimize for max watch time

transparency

- ✓
- ✓
- ✓
- ✓
- ✓

can be just serving customer needs/wants, not be relevant

Principal Components Analysis (PCA) - reduce # of features to be easier to combine (lose some info) visualize

features  $x_1, x_2, \dots, x_{100} \rightarrow$  new features  $z_1, z_2$

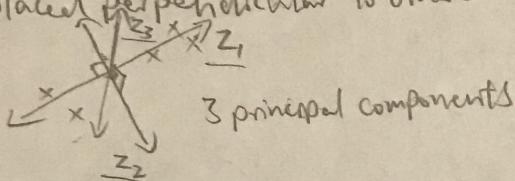
many features  $\rightarrow$  2 or 3 features  
# features and captures data

uses new  $z$ -axes values

① Normalize, feature scaling to data

② choose  $z = \alpha_0 + \beta_1 x_1 + \beta_2 x_2$  called principal component

each new principal component (axes) is placed perpendicular to others



- PCA is NOT linear regression

Linear Reg. minimizes distance along  $y$ -axis, whereas PCA treats all features  $x$  equally and reduce # of axes and maximize variance of points

"Reconstruction" - given  $z$ , reconstruct initial features

$$z = 3.55 \quad \text{original: } (2, 3)$$

$$3.55 \times \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix} = \begin{bmatrix} 2.52 \\ 2.52 \end{bmatrix} \rightarrow \text{estimates } (2.52, 2.52)$$

can't get original features, only can estimate

### Implementing PCA

① perform feature scaling

② "fit" data to obtain 2 or 3+ principal components

③ examine how much variance the principal components have (how much original data is preserved)

④ transform data onto new axes  $\rightarrow$  features will be reduced, can plot now

(old)  
PCA useful for visualization mainly, can be used for data compression or speeding up training models

Reinforcement Learning - training model to make decisions w/ max reward, given states ( $s$ )

$(s, a, R(s), s')$   $\rightarrow$  action of model

state action reward  $\uparrow$  new state

"stop"/done

of the current environment.  
for ambiguous contexts  
 $(+1, -10)$

$$\text{return} = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots \text{(until terminal state)}$$

Discount factor  $\gamma = 0.9$

Policy  $\pi(s) = \alpha^{\text{gamma}}$  function that maps from a state to actions

Goal of reinforcement learning is to find policy  $\pi$  that tells the model what action  $a$  to take given state  $s$ .

return depends on reward & discount factor depends on actions model takes

Markov Decision Process (MDP) - future only depends on current state, nothing else

	Mars Rover	Helicopter	Chess
states	6 states	position of helicopter	position on board
actions	$\leftarrow \rightarrow$	how to move control stick	possible move
rewards	100, 0, 40	+1, -1000	+1, 0, -1
discount factor $\gamma$	0.5	0.99	0.995
return	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3$
policy $\pi$	$\begin{matrix} 100 & \leftarrow & \leftarrow & \leftarrow & \rightarrow & 40 \end{matrix}$	Find $\pi(s)=a$	Find $\pi(s)=a$

### State-Action Value Function (Q Function)

$Q(s, a)$  = Return after taking action  $a$  from state  $s$  and taking optimal action after

The action that yields the highest return at a given state is the optimal action  
best return from state  $s$  is  $\max_a Q(s, a)$  and  $a$  is the best action  $Q^*$  = optimal Q function

### Bellman Equation

$s$  = current state

$R(s)$  = reward in

$a$  = current action

return current state

$s'$  = state after taking

from  $s$

action  $a$

$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$

$a'$  = action you take

in state  $s'$

$\gamma = 0.5$

$s=2$

$a=\rightarrow$

$s'=3$

$s=4$

$a=\leftarrow$

$s'=3$

$s=4$

Continuous State Spaces - oftentimes, an agent can have many continuous features describing its state, not just a set # of states

Ex. Continuous state features for...

car: x-position  
y-position  
θ-angle direction  
 $\dot{x}$  x-velocity  
 $\dot{y}$  y-velocity  
 $\dot{\theta}$  θ-velocity

helicopter: x-position  
y-position  
z-height  
roll left/right tilt  
pitch forward/backward  
yaw direction tilt  
rate of change  
roll pitch yaw

States usually have many features ranging many values "Continuous state MDP"

Example: Lunar Lander - land LL safely w/ reinforcement learning

Actions: -nothing  
-left thruster  
-mark thruster  
-right thruster

States: [x-position  
y-position  
x-velocity  
y-velocity  
θ-tilt  
θ-angular velocity  
l-left grounded  
r-right grounded]

Goal: learn policy  $\pi(s) = a$  and pick action  $a$  that maximizes return  $R_{\text{goal}}$

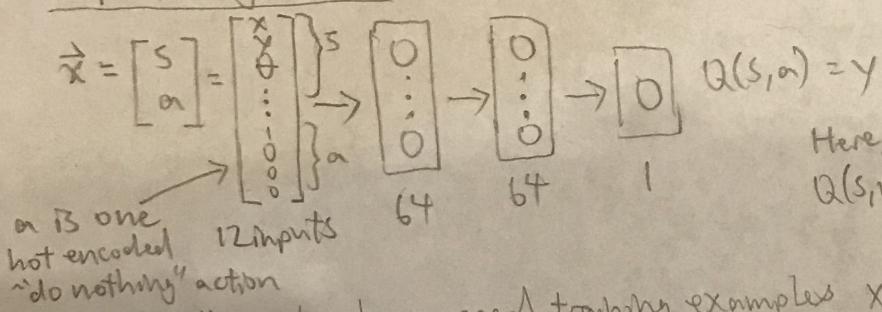
$\gamma = 0.985$

Rewards:

- Get to landing pad +100 - +140
- Reward for moving away/toward the pad
- Crash -100
- Soft landing +100
- Leg grounded +10
- Fire mark engine -0.3
- Fire side thruster -0.03

Calculating  $Q(s, a)$  - State-action value function (return)

Input state and action, output return  $Q_{\text{pred}}$  w/ neural net.



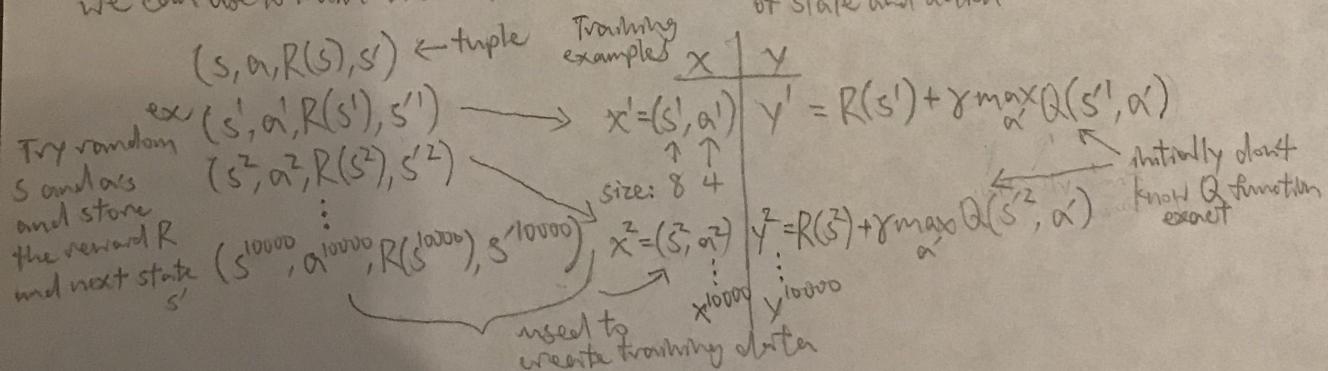
Here, compute all actions  $a$  at state  $s$ :  
 $Q(s, \text{nothing})$ ;  $Q(s, \text{left})$ ;  $Q(s, \text{mark})$ ;  $Q(s, \text{right})$

Pick action  $a$  that maximizes  $Q(s, a)$

To train the network, we need training examples  $x = (s, a)$  and labels  $y = Q(s, a)$

We use Bellman Equation to get data  $(x, y)$  that  $Q(s, a) = R(s) + \gamma \max_a Q(s', a)$

We can use to train neural network to predict return  $Q$  of state and action



Learning Algorithm to get Q-Function (Deep Q-Network (DQN)) - model that predicts return  $Q(s, a)$  for any  $s, a$

Initialize neural network with random weights for  $Q(s, a)$

Repeat, improving model:  $\sqrt{3}$ -Greedy Policy

Take random actions  $(s, a, R(s), s')$  and store recent 10,000 tuples  $\leftarrow$  Replay buffer

Train neural network: copy old weights to Direct NN, use mini-batch or "experience replay"

Create training set of 10,000 examples using Bellman & keep

$x = (s, a)$  and  $y = R(s) + \gamma \max_{a'} Q(s', a')$  or  $y = R(s)$  if terminal state  $\approx Q^*$  optimal  $Q$

Bellman Equation  $\leftarrow$  calculated w/ Direct

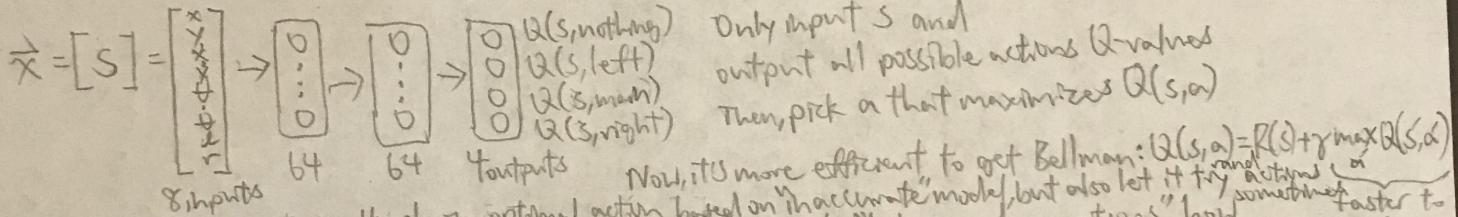
Train NN with dataset,  $Q_{\text{new}}$ , w/ new network w/ inaccurate params

$Q = Q_{\text{new}}$  w/  $\epsilon$ -Greedy policy  $MSE(y_{\text{new}}, Q)$  -  $Q$  will be inaccurate but improves from training NN initially over time

Use new, improved model, keep training it with slightly better model  $Q$   $Q = f_{\text{sub}}(x) \approx y$

### Architecture

NN Improvement - inefficient to run inference for  $Q$  for each action. Instead, simultaneously predict  $Q$  for all actions.



Now, it's more efficient to get Bellman:  $Q(s, a) = R(s) + \gamma \max_a Q(s, a)$

usually choose optimal action based on "inaccurate" model, but also let it try sometimes faster to get max

$\epsilon$ -Greedy Policy - When create training data, we said to "take random actions", how do we do that?

Because we randomly initialized parameters in NN, model initially predicts  $Q$  terribly and may

① with probability 0.95 choose action  $a$  maximizing  $Q(s, a)$  - Greedy, "Exploitation" not explore different actions.

② with probability 0.05 pick random action - "Exploration"

( $\epsilon=0.05$ ) - with probability 0.05 pick random action - "Exploration"

\* When model starts w/ terrible predictions, start  $\epsilon$  high. Gradually decrease  $\epsilon$

$$\epsilon = 1.0 \rightarrow 0.01$$

can be  
ex. 100,000,000

$\epsilon$  = probability to make random action

Algorithm Refinements: Mini-batch & Soft Updates

Mini-batch - regular gradient descent is slow b/c you sum all my examples for each step

- group examples in subset/batches "m", and each iteration of training, only use one subset  $\rightarrow$  also mini-batch

- converges slow but way fast computationally instead of using 10,000 examples, use 1,000 each iteration - faster! is used more often

Instead of using 10,000 examples, use 1,000 each iteration - faster!

Soft Update -  $Q_{\text{new}}$  (newly trained NN) can by chance be worse than  $Q$  (existing NN), so we don't want to completely overwrite it (maybe  $Q$  has some great weights)

$Q = Q_{\text{old}}$   $W = 0.01W_{\text{new}} + 0.99W$  no  $\rightarrow W = W_{\text{new}}$

keep some old params, change some to new  $B = 0.01B_{\text{new}} + 0.99B$  good  $\rightarrow B = B_{\text{new}}$

can set these #s Gradually converges to parameters for  $Q$  and is more reliable

## State of Reinforcement Learning

### Limitations

- way easier to work in simulation than real world (on pc vs physical robot)
- fewer applications than supervised and unsupervised learning
- is sick, research is being done, could have more applications in future

### Complete DQN Algorithm

create memory-buffer to store experiences, initialize  $\epsilon$ , # steps to update NN,

initialize<sup>2</sup> neural nets  $Q$  and  $Q_{new}$  w/ same random weights

repeat each full simulation "episodes" - times:

reset environment, keep initial state

for each time step, repeat:

pass state to  $Q$  NN to predict action, use  $\epsilon$ -Greedy Policy to choose action

take action in environment, store new state, reward, done

store experience tuple in memory-buffer

if condition to update (train) is met:

choose minibatch of <sup>random accumulated</sup> experiences (only every few steps, subset of experiences)

// train NN w/ both  $Q$  and  $Q_{new}$  NNs using minibatch data to use to train  $Q$  NN

adjust weights of  $Q$  NN w/ loss of Q prediction and  $Q_{new}$  = Bellman Equation,  $s = R(s)$  if terminal state,

$Q$  NN trained w/ minibatch gradient descent

use Soft Update to update  $Q$  NN weights with parts of weights of  $Q$  and  $Q_{new}$  or  $R(s) + \gamma \max Q_{new}(s', a')$  otherwise

after updating  $Q$  NN, set  $Q_{new}$  weights to  $Q$  to prep for next training

set state to next state

(keep repeating until episode ends)

update  $\epsilon$  value

Done!

## Topics Gone Over

### Supervised Machine Learning: Regression and Classification

Linear regression, logistic regression, cost function, gradient descent

### Advanced Learning Algorithms

Neural networks, decision trees, ML advice

### Unsupervised Learning, Recommenders, Reinforcement Learning

Clustering, anomaly detection, collaborative filtering, content-based filtering, reinforcement learning