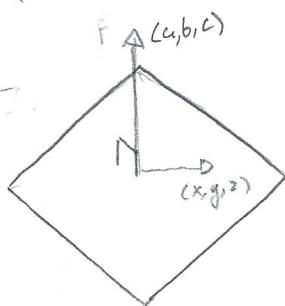
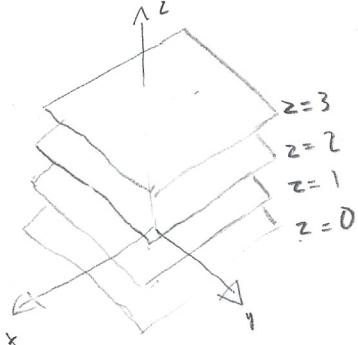


6.036 - Machine Learning.

Hollis Holmes.

Planes & Hyperplanes review:



Examples for intuition:

- equation $z=0$ defines the $x-y$ plane.
- if d is constant, $z=d$ is plane parallel to $x-y$ at $z=d$.
- equation $x+y+z=1$ defines slanted plane in \mathbb{R}^3 that goes through points $(1,0,0), (0,1,0), (0,0,1)$

Intercepts can be found by setting other coeffs to 0:
ex: z -intercept found $\Rightarrow ax+by+cz=d \Rightarrow z=\frac{d}{c}$.

Normal Vectors

If P is the plane in \mathbb{R}^3 defined by $ax+by+cz+d=0$ then:

$n = (a, b, c)$ is normal to this plane.

Distance From Point to Plane

- There is always 1 point on plane that is closest to the point. This point, q_0 , is the projection.
- Project the vector h (from any point on the plane, r , to point P) onto the normal vector from the plane.

$$|s| = |h \cdot u|, \text{ where } u \text{ is the unit normal vector}$$

Angle Between 2 planes (dihedral angle)

This is the angle between the plane's normal vectors.
This can be found from a dot product of the two vectors.

$$\text{Ex: } v = (1, 0, 1), w = (1, 1, 1)$$

$$v \cdot w = |v||w| \cos \theta \\ 2 = \sqrt{2} \sqrt{3} \cos \theta \Rightarrow \theta = \cos^{-1}\left(\frac{2}{\sqrt{2}\sqrt{3}}\right) \approx 35^\circ$$

Note for data representation & numpy:

- 2D arrays used to represent both matrices and vectors

Types:

Machine Learning - gave a data set with "answers", course ex data set w/ housing price for square footage

Supervised Learning - gave a data set with "answers", answer types: regression, classification,

Unsupervised Learning - data has no "right" answer, instead just looks for structure in the data, Market segmenting,

Models:

Linear regression: fitting predictive curve for real-value output.

$$\text{Size of house} \rightarrow h \rightarrow \text{price estimate}$$

$$h = \theta_0 + \theta_1 x$$

Notation: $M = \# \text{ training examples}$

$(x, y) = (\text{input, output})$

$x^{(i)} \rightarrow i^{\text{th}} \text{ training example.}$

Cost Function: how do we pick $h(x)$? create a cost function to quantify how bad the model

squared error example. $J(\theta_0, \theta_1) = \min_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$, square diff between output and estimate.

Gradient Descent: Algorithm to minimize the cost function and thereby, solve for $h(x)$.

walk in direction that has the most negative gradient until getting to minimum.

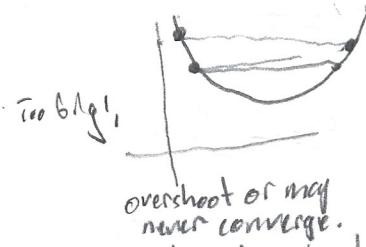
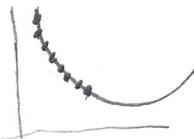
until at min:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

current change

what about alpha (α)? Too small:

too long to get there.

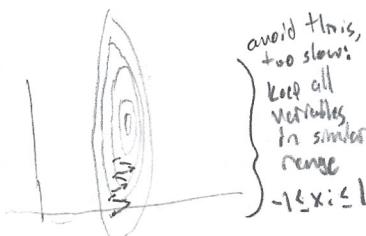


accidentally start at a local min?
grad. descent will have your value unchanged
good!

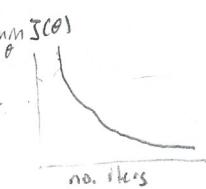
Week 2 As you approach a min, generally gradient will decrease, and you will automatically take smaller steps.

Multiple Features notation $x_j^{(i)}$ j^{th} feature of i^{th} example

Feature scaling: $x_1: \text{size}(0-2000)$
 $x_2: \text{bedrooms}(1-5)$
change to
 $x_1 = \text{size}/2000$
 $x_2 = \text{bedrooms}/5$.



Learning Rate (α , alpha) to help debug can plot cost against # iterations



$J(\theta)$ should decrease after every iteration.

Detect convergence if $J(\theta)$ decreases by less than ϵ at some iter
 $\epsilon: \epsilon = 10^{-3}$

Mean normalization: shift values to be centered about 0 ex: size $0-2000$
 $x_1 = \frac{x_1 - \bar{x}_1}{s_1}$ $\bar{x}_1 = \frac{1}{n} \sum x_1$

$$x_1 = \frac{x_1 - \bar{x}_1}{s_1}$$

$s_1 = \text{range}$

For sufficiently small α , $J(\theta)$ should decrease on every iter.
but if α is too small, grad. descent can be slow to converge.

Polynomial Regression Feature scaling becomes very important.

Normal Equations: Alternative to gradient descent, solves for optimal θ .
 $\frac{\partial}{\partial \theta_j} J(\theta) = 0 \quad \forall j$, all partial derivatives are 0 at minimum.
 Feature scaling not necessary.



(D) Grad. Descent

- Works well when n is large.

Normal Eqn

- No need to choose α
- No need to iterate.
- Have to compute $(X^T X)^{-1}$
 $\hookrightarrow n \times n \cdot O(n^3)$

$\theta = (X^T X)^{-1} X^T y$

$X^T X$ is rarely not invertible

\hookrightarrow pseudo inverse will still give right answer if not invertible.

Why not invertible: (singular, not full rank/col)

\hookrightarrow two features are related.

\hookrightarrow too many features ($m \leq n$)

\hookrightarrow malignant/tumor/benign.

Week 3.

Classification: {0, 1}

negative class positive class

regression isn't right to use here.

Logistic Regression

Want $0 < h(x) \leq 1$

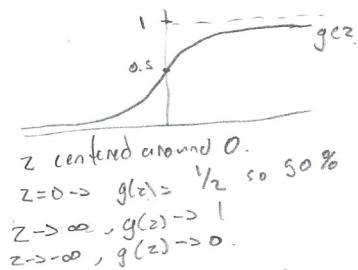
$$\text{sol: } h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Interpretation: $h_\theta(x)$ = estimated probability $y=1$ for this x .

$h_\theta(x) = P(y=1 | x; \theta)$ \Rightarrow probability $y=1$, given x , parameterized by θ .

Logistic Regression!

$$h_\theta(x) = g(\theta^T x), \quad g(z) = \frac{1}{1 + e^{-z}}$$



Cost function: want a convex cost function to guarantee finding a local min

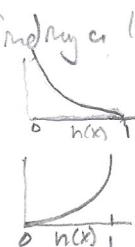
for logistic regression: $\text{Cost}(h(x), y) = \begin{cases} -\log(h(x)) & \text{if } y=1 \\ -\log(1-h(x)) & \text{if } y=0 \end{cases}$

- cost greater than 0
- increases further away you are from output.

Cost function: $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$

\hookrightarrow sum individual cost for all points.

$$\text{can be re-written: } J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \underbrace{y^{(i)} \log h_\theta(x^{(i)})}_{\text{if } y=1} + \underbrace{(1-y^{(i)}) \log (1-h_\theta(x^{(i)}))}_{\text{if } y=0} \right]$$

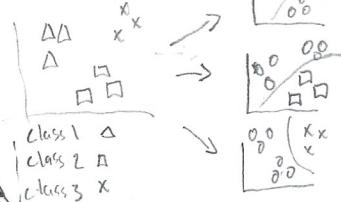


now take derivative for gradient descent:

$$\theta_j = \theta_j - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Multiclass Classification

: one vs all

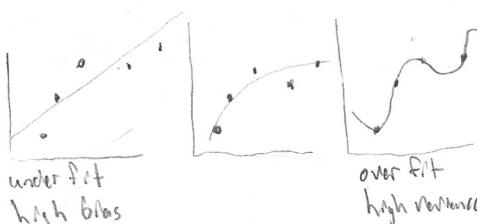


Solve a multiclass classification problem by solving k binary classification problems.

class i train logistic classifier $h_\theta^{(i)}(x)$
 for each class i to predict prob. $y=i$

then for new x , pick class s.t. $\max_i h_\theta^{(i)}(x)$

Overfitting



• overfit model has a very low cost function
 but fails to generalize to new examples.

Address overfitting: 1) Reduce # features. 2) Regularization, keep features but reduce magnitude/values of parameters θ_j

Regularization: penalizes making parameters θ_j big.

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

λ too high \rightarrow underfitting as parameters are too heavily penalized.

don't penalize θ_0

$$\text{logistic regression } J(\theta) = -\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)})) \right] + \lambda \sum_{j=1}^n \theta_j^2$$

Week 4

Neural Networks

- $a_i^{(l)}$ = activation of unit i in layer l
- $\theta^{(l)}$ = matrix of weights controlling function mapping from layer j to $j+1$

neuron

$$z_i^{(l)} \Rightarrow a_i^{(l)} = g(z_i^{(l)})$$

$$a_i^{(l)} = g(\theta_{10}^{(l)} x_0 + \theta_{11}^{(l)} x_1 + \theta_{12}^{(l)} x_2 + \theta_{13}^{(l)} x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(2)} x_0 + \theta_{31}^{(2)} x_1 + \theta_{32}^{(2)} x_2 + \theta_{33}^{(2)} x_3)$$

in previous layer + 1 $z_3^{(2)}$

$$\Rightarrow \theta^{(l)} = s_j \times (s_j + 1)$$

$$h_\theta(x) = a_1^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

Cost Function (classification problems) $L = \# \text{ layers}$, $s_k = \# \text{ units in layer } k$

$h_\theta(x) \in \mathbb{R}^k$, classifier problem

$$J(\theta) = \frac{-1}{m} \left[\sum_{i=1}^m \sum_{k=1}^{s_k} y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)}))_k \right]$$

$h_\theta(x)_j$ = j^{th} output.

- logistic regression cost over all training examples for the output node
- if multiclass, have to sum over all in the output layer.

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ij}^{(l)})^2$$

Don't sum over the bias (x_0) terms.

Backpropagation $\delta_j^{(l)}$: error of node j in layer l .

backprop calculates the direction to move for gradient descent, similar to what the partial derivatives do.

$$\delta_j^{(l)} = g_j^{(l)} - y_j^{(l)} \quad \text{have error term now work your way back}$$

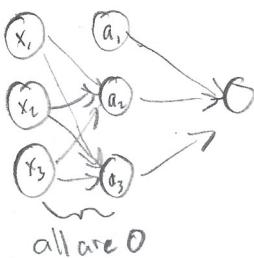
$$\delta_{j+1}^{(l+1)} = (\theta^{(l)})^T \delta^{(l)} \cdot g^{(l+1)}$$

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_i^{(l)} \delta_i^{(l+1)}$$

$$\frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \text{as an approximation to check gradient.}$$

Gradient Checking catch subtle errors in back propagation algorithm.

Initialising θ



since everything going from layer 1 \rightarrow layer 2 is the same (0)
 $\Rightarrow a_1^{(2)} = a_2^{(2)}$ and $\delta_1^{(2)} = \delta_2^{(2)}$ causes weights corresponding to the same input will all stay the same

To combat this Initialise all values in θ as random values in some range
 打破对称性 (good) also can't set all to the same non-zero value

Picking Architecture

- # of input units = dimension of features $x^{(i)}$
- # of output units = # of classes

reasonable default: 1 hidden layer or if $>$ hidden have same # of units in each hidden layer
 & usually the more the better.

Training A Neural Network: (may find local optima (non-convex))

- Randomly initialize weights
- Implement forward propagation to get $h_\theta(x^{(i)})$ for all $x^{(i)}$
- Implement code to compute cost function $J(\theta)$

- Implement backprop to compute partial derivatives $\frac{\partial}{\partial \theta_{jk}}$
- Use gradient checking to compare its estimate to $\frac{\partial}{\partial \theta_{jk}} J(\theta)$ from backprop
- Use gradient descent or also in backpropagation to try and minimize $J(\theta)$

1k 6 | Debugging / choosing a good implementation.

Imagine getting unreasonably large errors in your model's predictions

Machine learning diagnostic: test to gain insight into what is/ isn't working.

Evaluate A hypothesis predicted by model

split data set into a training set and test set on which trained model will be tested.

↳ Learn θ from training (minimize $J(\theta)$) then compute test set error.

Model Selection: need to pick degree of polynomial, "d".
So now we need to configure θ and "d".

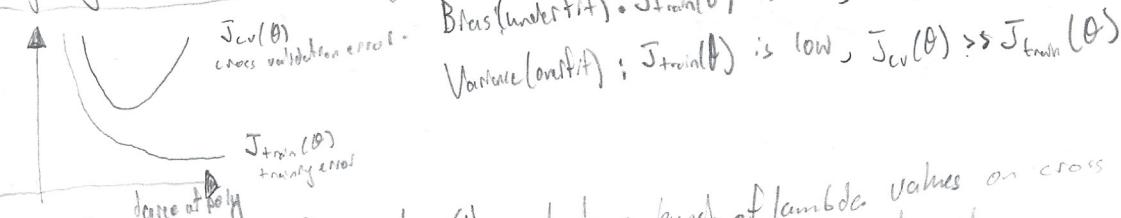
We calibrate θ on test set, so we need to test on a different data set to see if it generalizes well

So, split data into 3 parts (training, test, cross-validation)

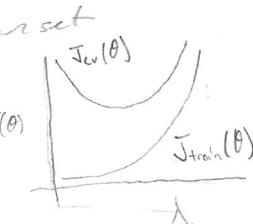
60% 20% 20%

- 1) Get set of θ 's for all degree polynomials running on training set
- 2) select d by checking error running each degree on cross-validation set
- 3) Test this model on test set

Diagnosing Bias Vs. Variance: if your model doesn't work as good as expected, you often have bias or variance issue



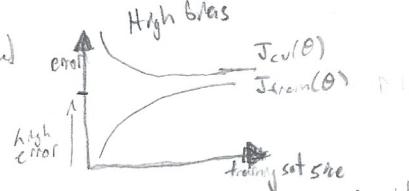
Choosing Regularization Parameter (λ), check a bunch of lambda values on cross validation set → train test on test set.



Learning Curves: good for sanity checks or performance of model

High bias: error does not improve w/ more data

High variance: getting more data can help.



Recommended Approach for new problem

- 1) start w/ very simple model that can be implemented and tested quickly → great for deciding what to work on next.
- 2) plot learning curves to decide next approach: more data, more features, etc
- 3) Error analysis: manually examine examples classified incorrectly (in cross validation set) to try and spot a pattern

Skewed Classes: Tumor classification, test set has 99% benign, 1% malignant tumors.

Boosting: query accuracy isn't good, because this could be accomplished by classifying everything as benign.

Need a new metric.

$$\text{Precision} = \frac{\text{True Positives}}{\text{Total Positives}}$$

$$\text{Recall} = \frac{\text{True Positives}}{\# \text{Actual Positives}}$$

Trading off Precision and Recall

Predict 1 if $h_0(x) \geq 0.3$ } high precision
Predict 0 if $h_0(x) \leq 0.3$

tag fewer wrong ones

Predict 1 if $h_0(x) \geq 0.3$ } high recall

Predict 0 if $h_0(x) \leq 0.3$

don't miss any.

Data in machine Learning: having more data can often improve model accuracy

ensure model has enough features to solve problem

Good sanity check is whether a human expert could do it with the information to the model.
Ex: Fill in this sentence, I have features, exp, housing price w/ only footage², no info about area or amenities.

of all patients who tested positive, what percentage are actually positive

of all patients that actually have cancer, what fraction did we correctly detect as having cancer.

	Precision	Recall
Alg01	0.5	0.4
Alg02	0.4	0.1
Alg03	0.02	1.0

How do we evaluate these?

↳ Average? $(P+R)/2 \rightarrow \text{BAD}$

↳ predict 0 or 1 all the time give recall or precision all the time but that is a useless classifier

Use the F-score: $\frac{P+R}{P+R}$

similar to average but gives the lower of precision and recall a lower weight

Hw 2 - Perceptron

1.1) If first a mistake $\rightarrow [1, -1]$, data 2 is correct, 3 is wrong so add it. $\rightarrow [-0.3, -2]$
 Now correct $\rightarrow 2$ mistakes.

b) $\{2, 3\}$

c) $\{0, -1\}$ then correct so 1 mistake. d) $\{1\}$

1.2) $a) x^{(3)} = [-10, -1], y^{(3)} = 1$
 $\{1, -1\}^{\textcircled{1}}, 2 \text{ good}, 3 \text{ wrong } \{ -9, -2 \}, 1 \text{ wrong } \{ -8, -3 \}, 2 \text{ right}, 3 \text{ right}, 1 \text{ wrong } \{ -7, -4 \},$
 $2 \text{ right}, 3 \text{ right}, 1 \text{ wrong } \{ -6, -5 \}, 2 \text{ right}, 3 \text{ right}, 1 \text{ wrong } \{ -5, -6 \}, \text{ all right}$
 $2 \text{ right}, 3 \text{ right}, 1 \text{ wrong } \{ -6, -5 \}^{\textcircled{2}}$
 $\rightarrow 6 \text{ mistakes}$

b) $\{0, -1\}, \text{ all right} \rightarrow 1 \text{ mistake}$

2) θ is initialized to the first mistake in the data set large $\theta \Rightarrow \text{large } \|x^{(i)}\| \Rightarrow \text{large } R$.

mistakes is a function of $R^2 \Rightarrow$ more mistakes

2.2) start w/ a correct answer difference result from 1 $\Rightarrow [1, -10]$

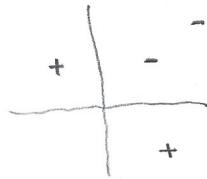
2.2) start w/ a correct answer difference result from 1 $\Rightarrow [1, -10]$

3) $\{2(-3) - 4(-1) - 2(-1) - 2(1), 2(2) - 4(1) - 2(-1) - 1(2)\} = \{-2, 0\} \quad \theta_0 = -5$

4) w/ 1 if $(1, 1, 1)$ b) not separable through origin. c) $\{1, 1, 1, -2.5\}$

{0 0.w
Do the types of classifiers work?

4.2)

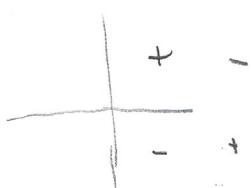


Do the types of classifiers work? (1.5, 1.5, 1)

- 1 circle \rightarrow no
- 2 circle centered anywhere \rightarrow yes \rightarrow circle around -ve points
- 3 line through origin \rightarrow no
- 4 line w/ offset $\rightarrow (-1, -1, 1)$

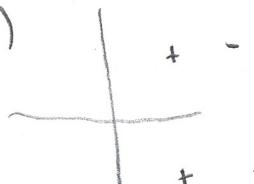
Which are lin-classifiers? $\rightarrow \{3, 4\}$

5.1)

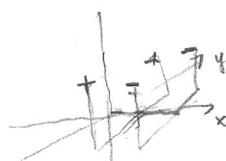


not lin-separable

5.2)



lin-separable, w/ offset



lin sep through origin

lin sep w/ offset

\rightarrow lin sep w/o offset

add extra dimension
to each point using
some non-zero #

6) R - bound on element magnitude, a) $\{1, 6.25\}$ b) $\{1, 0.5\}$ c) $\{1, 0\}$ d) $\{10, 2\}$
 γ - margin of dataset.

Mistakes bounded by $(\frac{R}{\gamma})^2$

Rest Coded

Week 3 - Features

Exercises.

1) One-hot rep

2)

lin sep $\bar{w} : |x|, x^4, x^{2k}$



Homework

1A) done in code

$$B) \left(\frac{R}{\delta}\right)^2 = 891,0670 \quad E) \left(\frac{1,2806}{0,0003}\right)^2 = 18222,229 \Rightarrow \text{more mistakes}$$

$$F) \left(\frac{R}{\delta}\right)^2 \Rightarrow R = 1,51 \Rightarrow \left(\frac{R}{\delta}\right)^2 \approx 32$$

4.2B) The most impactful parameters are cylinders (one-hot) and weight (standard). If two parameters could yield similar accuracy it would be these 2.

5.2) Reviews analysis. Need the 10 words with highest positive influence

→ Word that has highest value in θ

We need indices of θ that have highest value \rightarrow np.argmax()

5.2C) Found the highest and lowest reviews.

(optional) What is interesting is that they were very long. I think this is important. Though the sentiment of these reviews is generally positive/negative, the length of the reviews allow the margin to compile resulting in an extremely positive/negative review. It may be an improvement to account for the length of the review to get a better idea of something like the "density" of positives/negatives in the reviews. In this case I believe it is also important to remove stop-words, because these will unfairly decrease the magnitude of the sentiment in these longer reviews.

I later filtered the reviews for ones less than 500 characters long and got much more extreme reviews. This leads me to believe that this "density" property is important.

6.2F) Having a balanced data set is very important. For instance you could guess with 99% accuracy whether someone has cancer simply by guessing no for every person on Earth. However, this does not mean its a good test. The MNIST data set is balanced.

$$3) \begin{array}{c|ccccc} & + & + & + & + & x^2 \\ \hline -1 & & & & & \\ 0 & & & & & \\ 1 & & & & & \\ \hline & + & + & + & + & x \\ \end{array} \quad x^2 - 0.25 > 0 \quad |x| > 0.5 \\ -0.5, 0.5 \text{ both are the classified.}$$

Week 4 - Margin Maximization

Exercises: $\theta = (1, 1, -4)$: $(3, 2) \rightarrow \frac{1}{\| \theta \|}$, $(1, 1) \rightarrow \frac{2}{\| \theta \|}$, $(4, 2) \rightarrow \frac{-2}{\| \theta \|}$
 → Note θ_0 not included in $\| \theta \|$

Lab - Gradient Descent $x^{(k+1)} = x^{(k)} - \eta \nabla_x f(x^{(k)})$ $f(x) = (2x + 3)^2$

$$\text{local min } f'(x) = 0 \quad f'(x) = 2(2x + 3) \cdot 2 \Rightarrow 2x + 3 = 0 \Rightarrow x = -\frac{3}{2}$$

$$(1) \quad f'(x) = 8x + 12$$

written in these terms $\left. \begin{array}{l} x^{(k+1)} = x^{(k)} - \eta (8x^{(k)} + 12) \\ x^{(k+1)} + \frac{3}{2} = x^{(k)} - \eta (8x^{(k)} + 12) + \frac{3}{2} \\ = x^{(k)} (1 - 8\eta) + \left(\frac{3}{2} - 12\eta\right) \\ x^{(k+1)} + \frac{3}{2} = (1 - 8\eta)(x^{(k)} + \frac{3}{2}) \end{array} \right\} \quad \square$

$$z^{(k+1)} = \alpha z^{(k)}$$

where $z = x + \frac{3}{2}$
 and $\alpha = (1 - 8\eta)$

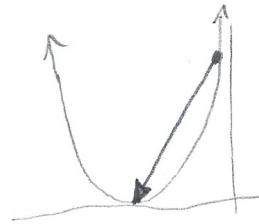
This tells us some interesting stuff about the starting point of our gradient descent:
 $|\alpha| > 1 \Rightarrow \alpha < -1 \text{ or } \alpha > 1$

$$|\alpha| > 1 \Rightarrow \text{diverges}; \eta < 0$$

this makes sense because we would be moving opposite the J'(x). calculated if we have $\eta < 0$

$\alpha = 1 \Rightarrow$ no gradient descent
 this makes sense no steps taken because $\eta = 0$
 $\alpha = -1 \Rightarrow$ oscillates endlessly between 2 numbers

$1 > \alpha \geq 0$ converges w/o oscillation
 $-1 < \alpha < 0$ converges w oscillation.



cool!

The maximum step size that converges w/o oscillation solves gradient descent in 1 step.

(2) $\eta = 0.1 \Rightarrow \alpha = 1 - 0.8 = 0.2 \Rightarrow$ converges w/o oscillation,

$$(3) \quad \eta = 0.1 \Rightarrow \alpha = 0.2$$

$$\eta = 0.11 \Rightarrow \alpha = 0.12$$

$$\eta = 0.12 \Rightarrow \alpha = 0.04 \quad \text{converges fast}$$

$$\eta = 0.13 \Rightarrow \alpha = -0.04$$

$$\eta = 0.14 \Rightarrow \alpha = -0.12 \quad \text{converges slow}$$

$$\eta = 0.15 \Rightarrow \alpha = -0.2$$

4(BS) hinge loss

↳ data isn't separable but still want to encourage margins.
 $L_h(\frac{\gamma(x, y, \theta, \theta_0)}{\delta_{ref}}) = \begin{cases} 1 - (\gamma(x, y, \theta, \theta_0) / \delta_{ref}) & \text{if } \gamma \geq \delta_{ref} \\ 0 & \text{else.} \end{cases}$ Loss is 0 if $\gamma > \delta_{ref}$ otherwise it is proportional to γ / δ_{ref}

Adding regularizer and letting $\gamma_{ref} = \frac{1}{\| \theta \|}$ we get objective function.

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n L_h(y^{(i)}(\theta^\top x^{(i)} + \theta_0)) + \lambda \| \theta \|^2$$

writing the training objective as an average:

$$\frac{1}{n} \sum_{i=1}^n \left[L_h(y^{(i)}(\theta \cdot x^{(i)})) + \frac{\lambda}{2} \| \theta \|^2 \right]$$

where

$$L_h(\theta \cdot x) = \max\{0, 1 - y(\theta \cdot x)\}$$

5(A) considering one point where hinge loss is positive this becomes: $J_\lambda(\theta) = 1 - y(\theta \cdot x) + \frac{\lambda}{2} \| \theta \|^2$

$$\text{define w.r.t } \theta \text{ and set to 0 for } \hat{\theta} \rightarrow \frac{\partial J}{\partial \theta} = -y x + \lambda \| \theta \| = 0 \Rightarrow \hat{\theta} = y x$$

□

Week 4 continued...

SB

We have: $\hat{\theta} = \frac{y^x}{x}$ when hinge loss is non-zero.
 general objective function for 1 point is: $J_\lambda(\theta) = [1 - y(\theta \cdot x)] + \frac{\lambda}{2} \|\theta\|^2$

$$\begin{aligned} \text{smallest } \hat{\theta} \text{ for which } L_h(y(\theta \cdot x)) = 0 \\ y(\hat{\theta} \cdot x) = 1 \Rightarrow \hat{\theta} \cdot x = \frac{1}{y} \quad \left| \begin{array}{l} \hat{\theta}^T x = \frac{1}{y} \\ \hat{\theta}^T x = \frac{x^T}{y} \end{array} \right. \Rightarrow \hat{\theta} = \frac{x^T}{\|x\|^2 y} \end{aligned}$$

SC) $\hat{\theta} = \hat{\theta}(\lambda)$ minimizes $J_\lambda(\theta)$

$$y(\hat{\theta} \cdot x) = y\left(\frac{y^x}{x}\right) \cdot x = \frac{\text{norm}(x)^2}{\lambda} \rightarrow \text{This term is always positive} \rightarrow \text{can't be misclassified.}$$

SD) single training element on the boundary, hence $1 - y(\theta \cdot x) = 0$
 point needs to be on the boundary

① hence: $y \cdot (\theta \cdot x) = 1$

we also need θ that minimizes our objective function, we solved for this $\hat{\theta}$

② hence: $\theta = \hat{\theta} = \frac{y^x}{x}$

$$\begin{aligned} \textcircled{1} \rightarrow \textcircled{2} \quad y\left[\frac{(y^x)}{\lambda} \cdot x\right] = 1 \\ y \cdot y = 1 \quad x \cdot x = \lambda \Rightarrow \lambda = \|x\|^2. \end{aligned}$$

M.2) SVM Gradient.

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \left[\begin{cases} 1 - y(\theta x_i + \theta_0) + \lambda \|\theta\|^2 & \\ 0 & \end{cases} \right]$$

$$\frac{\partial J(\theta, \theta_0)}{\partial \theta_0} = \theta_0 \left[\frac{1}{n} \left(-y x + 2\lambda \theta \right) \right]$$

Week 5 - Regression HWk, written work

1) $J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n L(x^{(i)}, y^{(i)}, \theta, \theta_0) + \lambda R(\theta, \theta_0)$

$$L_s(x^{(i)}, y^{(i)}, \theta, \theta_0) = (\theta^T x^{(i)} + \theta_0 - y^{(i)})^2$$

$\frac{\partial L_s}{\partial \theta} = 2(\theta^T x^{(i)} + \theta_0 - y^{(i)}) \cdot x^{(i)}$, $\frac{\partial L_s}{\partial \theta_0} = 2(\theta^T x^{(i)} + \theta_0 - y^{(i)})$
That was for individual points, now using X, Y which are the whole data sets.

$X: \delta x_n, Y: 1x_n, \theta: \delta x_n$

$$J(X, Y, \theta, \theta_0) = \underbrace{(\theta^T X + \theta_0 - Y)}_{\substack{1x1 \\ 1x_n}} \underbrace{(\theta^T X + \theta_0 - Y)^T}_{\substack{1x_n \\ 1x1}} \approx \text{LR} \rightarrow \frac{1}{n} \|(\theta^T X + \theta_0 - Y)\|^2$$

need to square & sum

$$\frac{\partial J}{\partial \theta} = \underbrace{X(\theta^T X + \theta_0 - Y)^T}_{\substack{\delta x_n \\ n \times 1}} + \underbrace{X^T(\theta^T X + \theta_0 - Y)}_{\substack{1 \times n}}$$

this is $1 \times d \rightarrow$ need dot
so: transpose the whole thing

- 2) Sources of error
- Structural: hypothesis class cannot represent the data (data too good)
 - Estimation: not enough data to accurately build a hypothesis. (model too good)

3) $\hat{\theta} = \frac{2}{n} \cdot Z^T \cdot (Z\theta - T)$

$$\hat{\theta} = Z^T Z \theta - Z^T T$$

$$Z^T Z \theta = Z^T T$$

$$\hat{\theta} = (Z^T Z)^{-1} Z^T T$$

$$\text{or } (X^T X)^{-1} X^T Y$$

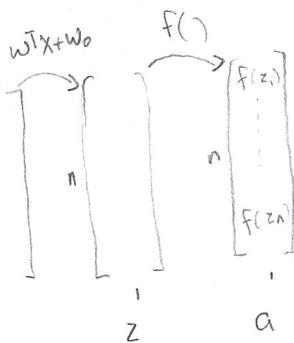
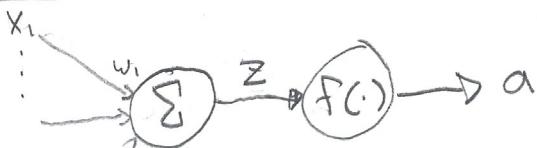
- 5) minimizing over
- dataset
 - predictor used
 - regularization?

Testing a data set don't need a regularizer, that is only for training a classifier.

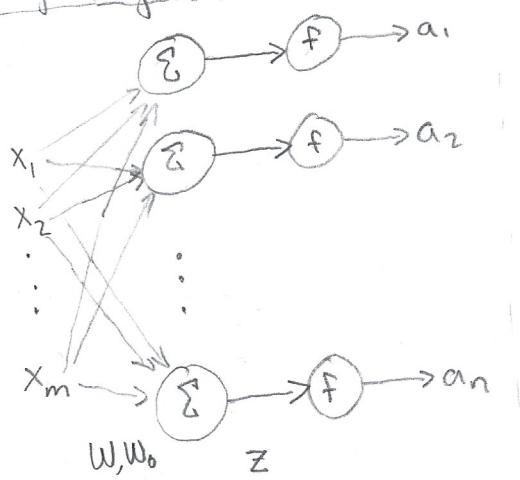
Week 6 - Neural Networks.

Basic Element

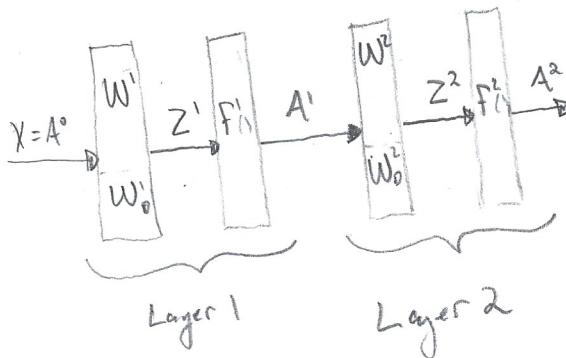
$$a = f(z) = f(\omega^T x + x_0)$$



Single Layer



Many layers



Back Propagation

In a neural net, the parameters are ω , so that is always what we need to calculate how much we need to change, and what we are calculating the gradient w.r.t. i.e: $\frac{\partial \text{Loss}}{\partial \omega}$. How much we need to change weights in order to reduce the loss incurred on this particular example.

First, looking at how loss depends on the final layer, using chain rule:

$$\frac{\partial \text{Loss}}{\partial \omega^L} = \underbrace{\frac{\partial \text{Loss}}{\partial A^L}}_{\text{derivative of Loss Function}} \cdot \underbrace{\frac{\partial A^L}{\partial Z^L}}_{f'(z)} \cdot \underbrace{\frac{\partial Z^L}{\partial \omega^L}}_{Z^L = (\omega^L)^T X^L} \Rightarrow \frac{\partial \text{Loss}}{\partial \omega^L} = A^{L-1} \left(\frac{\partial \text{Loss}}{\partial Z^L} \right)^T$$

$$\Rightarrow \frac{\partial Z^L}{\partial \omega} = X^L = A^{L-1}$$

Week 6 - Homework

$$1.1) z_i = w^T x_i$$

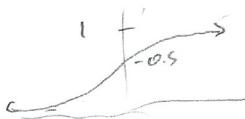
$$a = z_i$$

$$L = \begin{cases} 0 & \text{if } y_i a \geq 1 \\ 1 - y_i (w^T x_i) & \text{else} \end{cases}$$

$$\left. \begin{array}{l} L = \begin{cases} 0 & \text{if } w^T x_i (y_i) \geq 1 \\ 1 - y_i (w^T x_i) & \text{otherwise} \end{cases} \\ \frac{\partial L}{\partial w} = \begin{cases} 0 & \text{if } w^T x_i (y_i) \geq 1 \\ -y_i (x_i) & \text{otherwise} \end{cases} \end{array} \right. \quad w^T x_i (y_i) > 1$$

1.2) Log loss

$$\text{Sigmoid: } \sigma(z) = \frac{1}{1 + e^{-z}}$$



$$\text{ReLU: } \sigma(z) = \max(0, z)$$



$$\frac{\partial \text{Sigmoid}}{\partial z} = \frac{e^{-z}}{(1 + e^{-z})^2}$$

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

$$\frac{\partial (\text{sig})}{\partial z} = \text{sig}(1 - \text{sig})$$

D) For a single point. $NLL(a, y) = y \cdot \log(a) + (1-y) \cdot \log(1-a)$

$$\frac{\partial NLL}{\partial w} = \frac{\partial NLL}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w} =$$

$$NLL = -y \log(a) - (1-y) \log(1-a)$$

$$\frac{\partial NLL}{\partial w} = -y \frac{a(1-a)}{a} x_i f(1-y) \frac{a(1-a)}{1-a} x_i$$

$$= -y(1-a)x_i + (1-y)a x_i$$

$$= \{-y a + (1-y)\} x_i$$

$$= -(a-y) x_i$$

El vector:

2) Multiclass classification. $a = \text{SM}(z)$

Ssoftmax: takes n_L pre-activation values (z^L) returns a vector with each components probability where $\sum_j a_j^L = 1$, prob dist. over the categories

$$a_j^L = \frac{e^{z_j^L}}{\sum_{k=1}^{n_L} e^{z_k^L}} \quad (\text{your } e^{\text{pre-activation}} \text{ over } \sum(\text{all } e^{\text{pre-activation}}))$$

$$a) \text{SM}([-1, 0, 1]) = \frac{[e^{-1} \ e^0 \ e^1]}{e^{-1} + e^0 + e^1} \quad [0.09, 0.24, 0.76]$$

$$b) NLLM = - \sum_{j=1}^{n_L} y_i \ln a_j^L$$

sum along the length of the vectors multiplying matching component
only 1 should be there though

$$\frac{\partial NLLM}{\partial w_{kj}^L} = x_i (a_j^L - y_i)^\top$$

HW 6 continued...

3) Neural Networks

f_1 activation - ReLU = $\max(0, z)$

f_2 activation - Softmax

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \end{bmatrix} = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix}$$

$$\begin{bmatrix} w_{01}^1 & w_{02}^1 & w_{03}^1 & w_{04}^1 \\ w_{01}^2 & w_{02}^2 & w_{03}^2 & w_{04}^2 \end{bmatrix} = \begin{bmatrix} -1 & -1 & -1 & -1 \end{bmatrix}$$

$$\begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \\ w_{41}^2 & w_{42}^2 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 1 & -1 \\ 1 & -1 \\ 1 & -1 \end{bmatrix}, \quad \begin{bmatrix} w_{01}^2 \\ w_{02}^2 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$$

$$z^1 = \begin{bmatrix} z_1^1 \\ z_2^1 \\ z_3^1 \\ z_4^1 \\ z_5^1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 3 \\ 14 \end{bmatrix} + \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} 2 \\ 13 \\ -4 \\ -15 \end{bmatrix} \xrightarrow{\text{ReLU}} f(z^1) = \begin{bmatrix} 2 \\ 13 \\ 0 \\ 0 \end{bmatrix}$$

$$z^2 = \underbrace{w^T x + w_0}_z - \underbrace{a_2}_{\text{Softmax}}$$

$$z^2 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ -1 & -1 & -1 & -1 \end{bmatrix} \begin{bmatrix} 2 \\ 13 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 6 \\ 2 \end{bmatrix} = \begin{bmatrix} 15 \\ -13 \end{bmatrix}$$

$$a_2 = \text{softmax } c = \begin{bmatrix} e^{15} \\ e^{13} \\ e^{-12} \\ e^{-13} \end{bmatrix} \approx \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Week 6 Exercise Rough Work

Exercises

$$X = \begin{bmatrix} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 1 & 1 & 1 \end{bmatrix}$$

$n \times m$

W
 $\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$
 $2 \rightarrow z$ mapping
 x pre-activates
 coming from

$$W^T \text{ extended } X$$

$$\begin{bmatrix} 1 & 0 & -0.5 \\ -1 & 0 & 1.5 \end{bmatrix} \begin{bmatrix} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} -0.5 & 0.5 & 1.5 \\ 1.5 & 0.5 & -0.5 \end{bmatrix} \rightarrow f(z) = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \end{bmatrix} \left\{ \begin{array}{l} \text{AND Truth} \\ \text{table.} \end{array} \right.$$

$$y [0 \ 1 \ 0]$$

$$2) z = w_1 x_1 + w_2 x_2 + w_0 \quad y = -1 \quad \text{Training } L_h(v) = \max(0, 1-v)$$

$$= 1(1) + 1(2) + 1$$

$$= 4$$

$$f(z) = 1 \neq y$$

$$\frac{\partial L_{\text{Loss}}}{\partial w} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w}$$

$$\text{Loss} = 1 - v = 3 \quad \frac{\partial (1-y^a)}{\partial a} = -ya$$

$$1) \frac{\partial L_{\text{Loss}}}{\partial a} = -1$$

$$2) a = f(z) = a = z \Rightarrow \frac{\partial a}{\partial z} = 1$$

$$3) z = w^T x$$

$$= [1 \ 1 \ 1] \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \frac{\partial z}{\partial w} = x$$

$$\left(\frac{\partial z}{\partial w_1} = 1, \frac{\partial z}{\partial w_2} = 2, \frac{\partial z}{\partial w_0} = 1 \right)$$

$$\frac{\partial L_{\text{Loss}}}{\partial w_1} = (-1)(1)(1) = -1$$

$$\frac{\partial L_{\text{Loss}}}{\partial w_2} = (-1)(1)(2) = -2$$

$$\frac{\partial L_{\text{Loss}}}{\partial w_0} = -1$$

$$w = [1 \ 1 \ 1] - 0.5[-1 \ 1 \ 2] = [0.5 \ 0.5 \ 0]$$

New output -2

$$\text{New loss } \max(0, 1 - (-1)(-2)) = 0$$

→ no new updates

$$w = [0.5, 0, 0.5] - 0.5[1, 2, 1]$$

$$= [0, -1, 0]$$

Week 6 Lab - Rough Work

1.1(A) $Z = \mathbf{w}^T \mathbf{x}$
 $a = \frac{1}{1 + e^{-Z}}$

separator is $a = 0.5$

$$0.5 = \frac{1}{1 + e^{-Z}}$$

$$0.5(1 + e^{-Z}) = 1$$

$$e^{-Z} = 1 \Rightarrow Z = 0$$

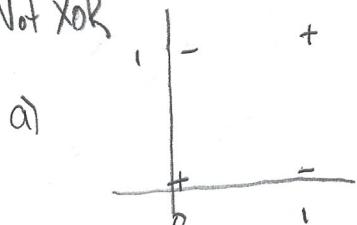
(Ans: Linear)

$$\mathbf{o} = \mathbf{w}^T \mathbf{x}$$

↳ This is a linear system of equations, hence the operator is linear.

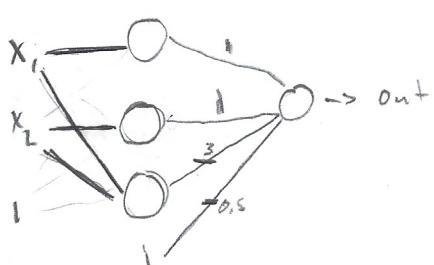
B) yes, increasing the iters by a factor of 10 yields consistent results.

1.2) Not XOR



No, the network above would not work as its hypothesis class is linear functions, but this data-set is not linearly separable.

b)



A layer w 3 hidden units can be used to build a classifier as described below

top unit represents if x_1 is a one.

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

the weights coming out of these require at least 1 of the top 2 to be on, but is 0 if both are on with a large negative

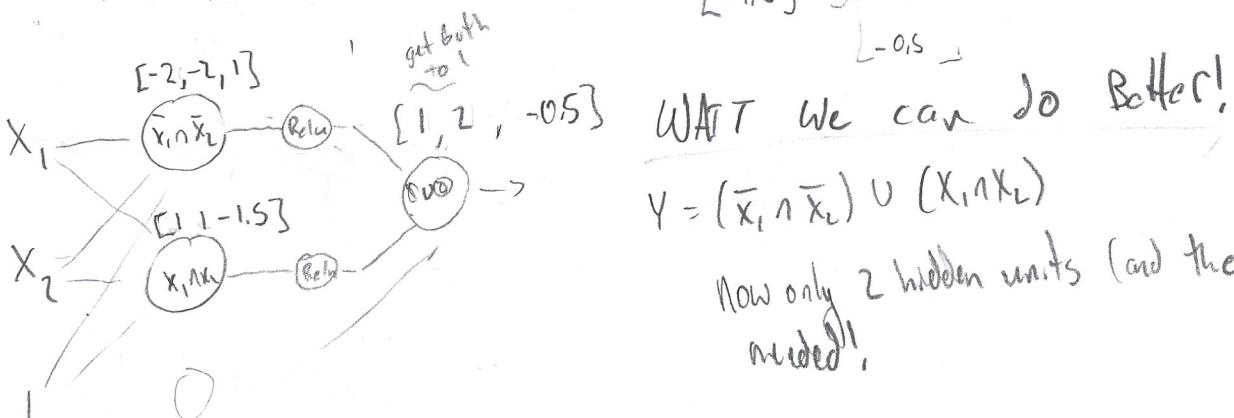
middle represents if x_2 is a one

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

bottom represents if x_1 and x_2 are one

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

$$-0.5$$



WAIT we can do better!

$$Y = (\bar{x}_1 \wedge \bar{x}_2) \vee (x_1 \wedge x_2)$$

Now only 2 hidden units (and the x0) are needed!

$$Y = (\bar{x}_1 \wedge \bar{x}_2) \wedge (x_1 \wedge x_2)$$

represent each of these in the hidden layer, then in the last.

Week 6 Lab Ctd'ns

13) Hard data set.

I was able to get a reliable 95% accuracy with parameters $hu=18$, $iters=5000$, $lr=0.001$. However, this value set was tough to find because of the small how close some of the +ve and -ve points are.

b) No I don't think its a good idea to find a perfect separator. By only maximizing training accuracy you will likely over-fit the data and the model will perform poorly on other data sets.

2) ① # output units, activation function on output units, and loss function
↳ should all be chosen jointly.

A) Words on front page NY Times \rightarrow change in stock market.

① one output unit, magnitude or percent change in price.

② Relu - bunch of parameters, not probabilistic

③ squared loss, care about the magnitude difference between prediction and a

B) Satellite image \rightarrow probability it will rain in next 4 hours.

① one output unit, probability

② sigmoid, map to a probability, binary classification.

③ Neg log. likelihood - maximize probability over the data set.

C) Words in an email \rightarrow folder it should go in.

① # of folders that can be chosen, one-hot classification.

② softmax - multiclass classification

③ NLL - maximize probability over the whole dataset

D) Words of a document \rightarrow vector of topics, each index is 1 if topic is addressed

① not strictly classification because vector can have multiple 1's.

② If there are T topics, this seems like T-copies of a binary classified since each topic is either mentioned or not-mentioned

③

④

⑤

HW6 cont'd...

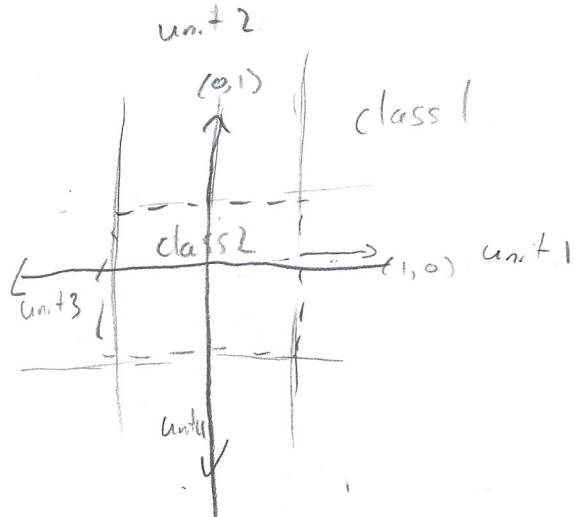
3.2) Unit decision boundaries.

$$\{1, 0\}[-1] \text{ unit 1 } x_1 - 1 = 0$$

$$\{0, 1\}[+1] \text{ unit 2 } x_2 - 1 = 0$$

$$\{-1, 0\}[+1] \text{ unit 3 } -x_1 - 1 = 0$$

$$\{0, -1\}[-1] \text{ unit 4 } -x_2 - 1 = 0$$



The intuition here is that each row in W^T is its own linear separator.

Now outputs of the whole system, the classifiers.

$$f'(z_1) + f'(z_2) + f'(z_3) + f'(z_4) = 0$$

$$z_1^2 = f(z_1) + f(z_2) + f(z_3) + f(z_4) + 0 = 0 \quad \left. \begin{array}{l} \\ z_1 = 2 \\ z_2 = -1 \end{array} \right\}$$

$$a_1^2 = \frac{e^0}{e^0 + e^2} = \frac{1}{2}$$

$$a_2^2 = \frac{e^2}{e^0 + e^2} = \frac{1}{2}$$

$$f'(z_1) + f'(z_2) + f'(z_3) + f'(z_4) = 1 \quad \boxed{\text{Boundary}}$$

$$z_1^2 = 1 + 0 = 1 \quad a_1^2 = \frac{e^1}{e^1 + e^1} = \frac{1}{2}$$

$$z_2^2 = -1 + 2 = 1 \quad a_2^2 = \frac{e^1}{e^1 + e^1} = \frac{1}{2}$$

$$f'(z_1) + f'(z_2) + f'(z_3) + f'(z_4) = 3 \quad \boxed{\text{Class}}$$

$$z_1^2 = 3 + 0 = 3 \quad a_1^2 = \frac{e^3}{e^3 + e^1} \approx 0.98$$

$$z_2^2 = -3 + 2 = -1 \quad a_2^2 = \frac{e^{-1}}{e^3 + e^{-1}} \approx 0.02$$

class 1 seems to just be are you outside of the 1×1 .

class 2 is are you pretty close to the origin.

In this case ReLU's are all 0 so point is in none of the unit zones
→ it is in the middle 1×1

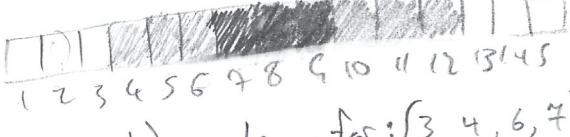
over one of the boundaries by or just over 2 of them sum to 1.
on the boundary equal prob of both classes

Way out in one of the zones, far away from the origin. \tilde{W} high prob

→ cool how a point can stay stationary but with only magnitudes of weights changing probability/confidence of model prediction changes.

Lab 8 - Convolutional NN's

2A) $x = \{0, 0, 0, 2, 2, 4, 4, 4, 3, 2, 2, 0, 0\}$

as gray scale 

2B $x_3 - x_1 > 1$ $(-1, 0, 1, -1)$ true for: $\{3, 4, 6, 7\}$
 $x_1 - x_3 > 1$ $(1, 0, -1, -1)$ $\{7, 8, 10, 11\}$

3A) $f_1 \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$; $f_2 \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$; $f_3 \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$

} Decreasing the 6.
 (making it more:
 means the feature
 being detected
 needs to be
 stronger for it
 to be detected

-detects vertical -detects horizontal -detects steeper pixels that are
 lines. lines. brighter than
 its surroundings.

D) changing coef's on the array other than bias.

If bias is 0: this has no effect

If bias is non-0: increasing coefficient reduces the affect of bias

Decreasing coef increases " " " " " "

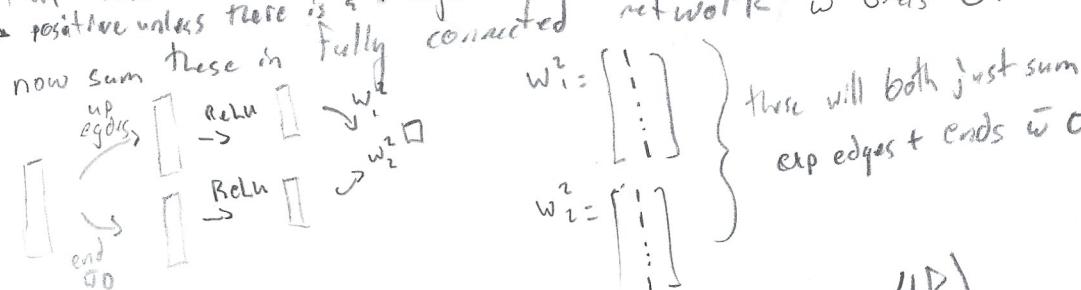
ex setting coef to 0 \rightarrow all pixels are just determined by bias. $0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0$

E) series of 0's is an obstacle. how do we detect this?

every change up 0 \rightarrow 1 means you came from a 0. $(-1, 1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1)$ -0.5

you then need to check if it ends in a 0. $(-1, -1, -1, \dots, -1, 0)$ 0.5

\rightarrow positive unless there is a 1 anywhere.



4) $\boxed{0} \boxed{x_1} \boxed{x_2} \boxed{x_3} \boxed{x_4} \boxed{0}$

Notice how parameters 'slide' and are reused multiple times in the matrix. This is parameters.

$$\begin{array}{l} 6x_1 + x_2 - 0x_3 + 0x_4 = z_1 \\ 5x_1 + 6x_2 + 7x_3 + 0x_4 = z_2 \\ 0x_1 + 5x_2 + 6x_3 + 4x_4 = z_3 \\ 0x_1 + 0x_2 + 5x_3 + 6x_4 = z_4 \end{array}$$

4B)

$$d - (k-1) + 2 \cdot p$$

Week 9 Lab

1) State space: {dirty, clean, painted}
 Action space: {wash, paint, eject}

$$T(s_t, \text{wash}, s_{t+1}) = \begin{bmatrix} 0.1 & 0.9 & 0 \\ 0.1 & 0.9 & 0 \\ 0.1 & 0.9 & 0 \end{bmatrix}, \quad T(s_t, \text{paint}, s_{t+1}) = \begin{bmatrix} 1 & 0 & 0 \\ 0.1 & 0.1 & 0.8 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R(s, a) = \begin{cases} 10 & \text{if } (s = \text{painted}, a = \text{eject}) \\ 0 & \text{if } (s \neq \text{painted}, a = \text{eject}) \\ -3 & \text{if } (a \neq \text{eject}) \end{cases}$$

b) if dirty \rightarrow wash, if clean \rightarrow paint, if painted \rightarrow eject

c) now you only survive with probability 0.1, so rewards now are heavily valued.

All plates start dirty so cleaning and painting earns:

$$-(3) - (0.1)(3) + (0.1)^2 10 = -3.2.$$

\rightarrow Doesn't even make sense to paint plates, just eject immediately for 0 reward.

d) If horizon is 2 can never eject a clean plate so just eject the dirty ones.

$$\begin{bmatrix} * \\ \$ \end{bmatrix}$$

2) A) \$ is 100, * is 50

B) states find each reward square. \$ is worth more given its higher one-time reward, and states haven't had a chance to redeem * multiple times.

c) States that were pointing towards * have now been able to reach \$ but not able to get back to * twice. So the $T*$ for these squares is shifting back towards \$.

d) Those states are now claiming * multiple times and it's shifting to going towards \$.

e) Almost all values have been able to claim * multiple times.

f) All values will point towards * and their values will keep climbing.

$$3) y_t = (1)y_{t+1} - (2)y_{t+2} + (3)y_{t+3}$$

$$s_t = f_1(W^s s_{t+1} + W^x x_t + W^o o_t)$$

$$y_t = f_2(W^s s_t + W^o o_t)$$

Week 10 Lab

- 2A) With Q-learning, the 'learning' doesn't start until the robots first hits the target square. Until then, all squares have 0 value (depending on implementation). With Value iteration, target states are given a value on the first iteration and the values propagate outwards.
- B) The robot has not come across the target-state yet. So, the value of all values is 0.
- C) Once square is yellow because the goal state has been found. After it has been found, robots position is re-initialized to a random point on the grid. No more updates are made until the robot stumbles upon a square adjacent to the target square.
- D) Once there are a few squares w/ values, the robot stumbles on these positions more and more quickly. This in turn, allows q-values to propagate out more quickly.
- E) Choosing a random action instead of the first in case of a tie.
- F) Directions are random and value of all squares are plummeting. Because each update the value gets reduced in the SSP model.
- G) As opposed to target-oriented, almost all values on the grid has value changing very often. This is because in goal-oriented only 1 square has a reward, but in SSP, all but 1 have a reward.
- H) I think this depends on your use. Goal-oriented can take a very long time to get started, but ended with a more complete sol than SSP ^{in this example}.

3) This is a tradeoff between exploration vs. exploitation.
Exploitation acts exactly in accordance with the current Q-table.
However, you could imagine that if you haven't spent enough
time exploring you will concentrate on a poor solution. Conversely,
at some point you would like to shift your efforts
towards focusing on areas you found to be important.

Epsilon search:

- Random choice in probability ϵ (exploration)
- Best according to current Q in prob $(1-\epsilon)$ (exploit)

With $\epsilon = 0$: purely exploitation, risks never finding the optimal policy, puts too much faith in results w/ very little data

With $\epsilon = 1$: purely exploration will find a good policy but never use it because it chooses at random always,

$\epsilon = 0.5$: a mix of both, some exploration, some exploitation.

None of these guarantee optimal behaviour during learning.