

COMP 540 Final Exam Review

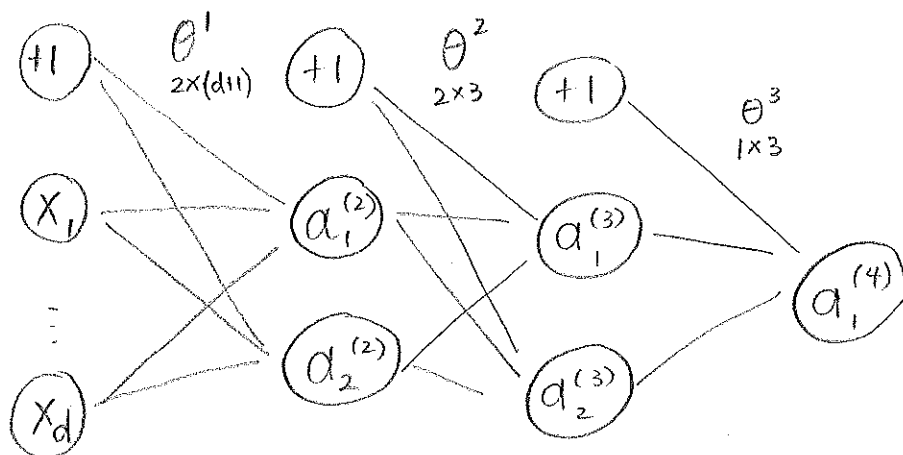
April 26

Basic Neural Networks

notation: $a_j^{(l)}$ - activation of j^{th} unit in l^{th} layer

$\theta^{(l)}$ - parameters mapping layer l to $l+1$

* the dimension of θ is always $\begin{matrix} \text{\#nodes in next layer (no bias)} \\ \times \\ \text{\#nodes in previous layer (includes bias)} \end{matrix}$



Forward Propagation

$$a^{(1)} = x \quad z^{(2)} = \theta^{(1)} * [1; x]$$

$2 \times (d+1) \quad (d+1) \times 1$

$$a^{(2)} = \text{ReLU}(z^{(2)})$$

$$z^{(3)} = \theta^{(2)} * [1; a^{(2)}]$$

$2 \times 3 \quad 3 \times 1$

$$a^{(3)} = \text{ReLU}(z^{(3)})$$

general step

$$z^{(l+1)} = \theta^{(l)} * [1; a^{(l)}]$$

$$a^{(l+1)} = \text{ReLU}(z^{(l+1)})$$

$$a^{(4)} = g(\theta^3 [1; g(\theta^2 [1; g(\theta^1 [1; x])])]) = g(h_\theta x)$$

* any \wedge continuous function can be approximated by 3 layers. 4 layers approximates any fcn

ost

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k \log(h_\theta(x))_k + (1 - y_k) \log(1 - h_\theta(x))_k + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{k=1}^{S_l} \sum_{j=1}^{S_{l+1}} \theta_{ji}^{(l)2}$$

$$\frac{\partial J(\theta)}{\partial \theta} = -\frac{1}{m} \sum_{i=1}^m \underbrace{(y^{(i)} - h_\theta(x^{(i)}))}_{\text{error in output}} x^{(i)}$$

K - # output units (final output)
 S_l - # units in layer l

Backward Propagation

Need: $J(\theta)$ and $\frac{\partial J}{\partial \theta_{ij}}$

addendum: for $i=1$ to m (training examples)

① set $a^{(1)} = x^{(1)}$

① compute forward prop $a^{(2)} \dots a^{(L)}$

② δ = error of node j in layer l , error in activation

ex: $\delta_j^{(4)} = a_j^{(4)} - y_j$, $\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} \cdot \underbrace{g'(z^{(3)})}_{\text{derivative}}$
 $a^{(3)} \cdot (1 - a^{(3)})$

- compute $\delta^{(L-1)} \dots \delta^{(2)}$ using $\delta^{(l)}$

③ $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

because of the fact that $\frac{\partial J}{\partial \theta_{ij}^{(l)}} = a_j^{(l)} \delta_i^{(l+1)}$ *delta vals are derivative of cost fn

(this is w/o regularization)

vectorized:

$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

end loop

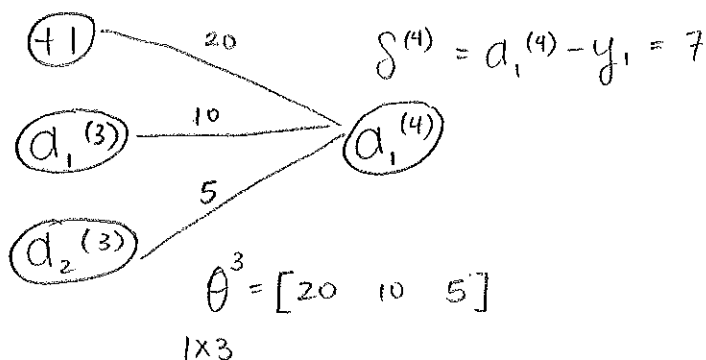
$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}$ if $j \neq 0$

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$ if $j = 0$

} accumulates partial derivatives

$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = D_{ij}^{(l)}$

Visually:



$\delta_1^{(3)} = \theta^{(3)T} \delta^{(4)} \cdot \underbrace{g'(z^{(3)})}_{\text{derivative}}$

$\delta_1^{(3)} = \begin{bmatrix} 20 \\ 10 \\ 5 \end{bmatrix}_{3 \times 1}^T \cdot \begin{bmatrix} 7 \end{bmatrix}_{1 \times 1} \cdot \begin{bmatrix} \times \\ \times \\ \times \end{bmatrix}_{3 \times 1}$

discard δ_0 for bias unit!!

NN Tactics, Specifications, Hyperparameters

Constructing a network =

- # of input units = dimensions of $x^{(n)} \in 1$ to d
- # of output units = # classes
- # of hidden units/layer
- recommend same # units in each hidden layer

Training a network:

1. randomly initialize weights
2. forward prop
3. implement cost function
4. back prop to compute partial derivatives
5. grad checking
6. gradient descent to min cost function

— ? unclear about steps

* for $i=1$ to m (do for each example)

FP
BP

Tips:

- mean-subtract, zero center, normalize (div by std)
 - get mean/std from train set and USE THEM for val/test
- Batch Normalization: normalize activation of units to have $\sim N(0,1)$ dist at beginning of training $a_j^{(l)} \sim N(0,1)$
 - reduces dependence of gradients on scale of parameters or init values
 - can regularize
 - can lead to faster convergence
- L2 reg
 - sum of squared errors
- L1 reg
 - sparse weight vectors
 - used for feature selection
- overfit on small subset of data
100% on small train
- monitor learning rate
- dropout = unit active w/ prob p
- dropconnect: some weights set to 0 during forward pass
- gradient checks!!
numerical approximations
$$\frac{f(x+h) - f(x-h)}{2h}$$

Update momentums (SGD, Nesterov, ADAM, etc)

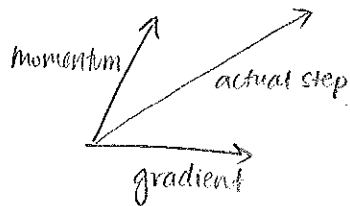
SGD $X += -\alpha \frac{\partial J}{\partial \theta}$

α = learning rate

Momentum $\mu \rightarrow v - \alpha \frac{\partial J}{\partial \theta}$

$$X += \mu * v - \alpha \frac{\partial J}{\partial \theta}$$

get over bumps in loss functions

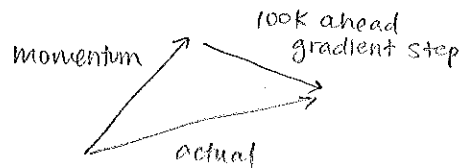


Nesterov momentum

$$X_{\text{ahead}} = X + \mu v$$

$$v = \mu v - \alpha \frac{\partial J}{\partial \theta} (\text{ahead})$$

$$X += v$$



Tips (more):

- annealing the learning rate
 - reduce step by 0.1 every 20 epoch
 - exponential decay

Adaptive Learning Rates

- 1) RMSProp: moving avg of squared gradients
- 2) Adam: smoothed RMS Prop w/ momentum
- 3) Adagrad: also depends on gradients

Convolutional Neural Networks

- a network learns filters that are local in width/height but full in depth
- hyperparams: depth, stride, and zero-padding

- 1) depth: # units in CONV that connects to one same input region
- 2) stride: small stride = largely overlapping receptive fields / large output
- 3) padding: pad input w/ zeros all around \rightarrow control spatial size

CONV

- input volume of size $W1 \times H1 \times D1$

- # filters K
- receptive field size F
- stride S
- zero-pad P

- produces output volume of size $W2 \times H2 \times D2$

$$W2 = \left\lfloor \frac{W1 - F + 2P}{S} \right\rfloor + 1$$

$$H2 = \left\lfloor \frac{H1 - F + 2P}{S} \right\rfloor + 1$$

$$D2 = K$$

$F \times F \times D1$ weights per filter

total $F \times F \times D1 \times K$ weights and K bias units

* can dramatically reduce # params by making assumption that the same weights and biases are used across depth slices

Common to insert pool layer after conv.

- 2×2 filter w/ size 2 \rightarrow take max over 4 numbers
- depth is same

params: • receptive field size F
• stride S

pooling is
downsampling!!

POOL

input vol: $W1 \times H1 \times D1$

output: $W2 \times H2 \times D2$

$$W2 = \frac{(W1 - F)}{S} + 1$$

$$H2 = \frac{(H1 - F)}{S} + 1$$

$$D2 = D1$$

More Notes about NN

activation functions:

1) sigmoid (σ)

- saturate & kill gradients at 0 or 1
- not zero centered outputs (they are 0.5)

2) tanh

- also saturate activations
- squash $[-1, 1]$
- zero-centered !! ✓

3) ReLU

- $f(x) = \max(0, x)$
- linear, non saturating
- accelerates SGD convergence
- gradients could die b/c of gradients that are huge \rightarrow push to zero

4) Leaky ReLU

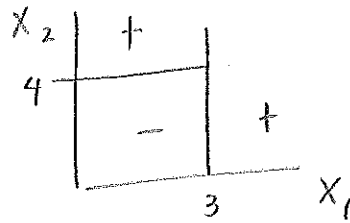
- ReLU w/ a $x < 0$ slope component

Conv Net Tips

- Data augmentation
- ReLU
- dropout
- mini-batch grad descent
- take out layers and see what impact it has on accuracy
- too deep: weights at beginning become meaningless + small

April 27

ex:



Decision Trees

- type of adaptive basis function

$$h_\theta(x) = \theta^T \phi(x) \quad \text{where} \quad \phi_1(x) = I[x_1 < 3 \text{ and } x_2 > 4] = +$$

$$\phi_2(x) = I[x_1 < 3 \text{ and } x_2 \leq 4] = -$$

$$\phi_3(x) = I[x \geq 3] = +$$

$$h_\theta(x) = \theta_1 \phi_1(x) + \theta_2 \phi_2(x) + \theta_3 \phi_3(x)$$

- Must pick attribute that allows greatest reduction in cost
- too deep of a tree \rightarrow overfitting
- highly unstable
- depth = hyperparameter. use validation to prune on overfit model

x_j^* 's - features

t^* 's - thresholds

feature selection: (maximizing reduction in cost due to split)

$$j^*, t^* = \arg \max_{\substack{j \in \{1, \dots, d\} \\ t \in \text{values of } x_j}}$$

$$\left\{ \text{cost}(\mathcal{D}) - \left[\frac{|\mathcal{D}_{\text{left}}|}{|\mathcal{D}|} \text{cost}(\mathcal{D}_{\text{left}}) + \frac{|\mathcal{D}_{\text{right}}|}{|\mathcal{D}|} \text{cost}(\mathcal{D}_{\text{right}}) \right] \right\}$$

cost of split

reduction in cost due to split = info gain due to split

Regression cost $y \in \mathbb{R}$

$$\text{cost}(\mathcal{D}) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \bar{y})^2$$

every feature = $O(D)$ in computation

classification cost

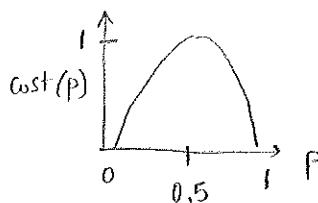
$$y \in \{0, 1\}$$

$$\text{cost}(\mathcal{D}) = \frac{1}{m} \sum_{i=1}^m I(y^{(i)} \neq \hat{y})$$

$$\hat{y} = \arg \max_c I(y^{(i)} = c)$$

Entropy

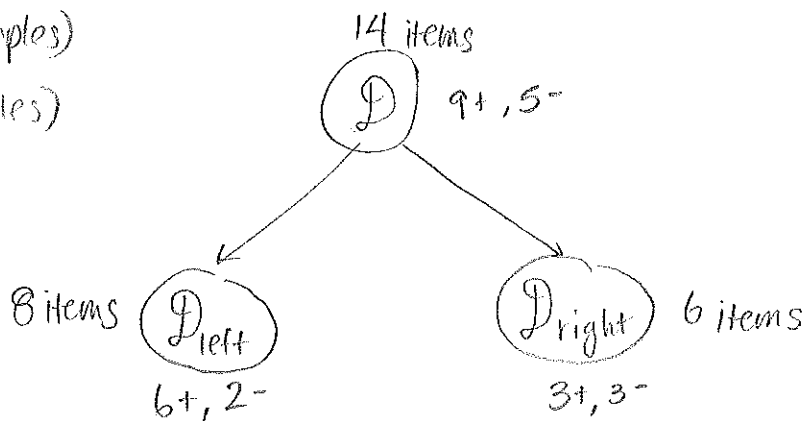
$$\text{cost}(D) = -p \log p - (1-p) \log (1-p)$$



$p = \# \text{ pos examples}$

example: 9 + (pos examples)

5 - (neg examples)



$$\text{cost}(D) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

$$\text{cost}(D_{\text{left}}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0.811$$

$$\text{cost}(D_{\text{right}}) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$$

$$\text{Information gain: } \text{cost}(D) - \left[\frac{|D_{\text{left}}|}{|D|} \text{cost}(D_{\text{left}}) + \frac{|D_{\text{right}}|}{|D|} \text{cost}(D_{\text{right}}) \right]$$

$$= 0.94 - \underbrace{\left[\frac{8}{14} (0.811) + \frac{6}{14} (1) \right]}_{\text{cost of split}} = \text{reduction in cost OR info gain}$$

Gini Index

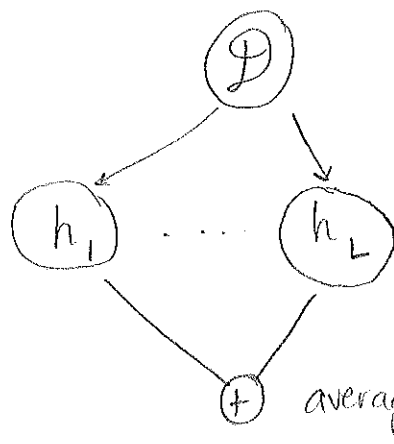
$$\text{cost}(D) = 2p(1-p)$$

$$p = \frac{9}{14}$$

* read more on this!!

$$= \frac{18}{14} \left(\frac{5}{14} \right)$$

Ensemble Models



bootstrapping/bagging
 \Downarrow
 Sample w/ replacement
 m times

$$h_{\text{bag}}(x) = \frac{1}{L} \sum_{l=1}^L h_l(x)$$

$$h_l(x) = \theta_l^T x$$

$h_{\text{bag}}(x) \rightarrow$ majority vote among $h_1(x) \dots h_L(x)$

Serves to uncorrelate
 samples and errors
 (not always perfect
 in reality)

* Error of bagged ensemble is lower than average expected error

classification:

$$P(\text{ensemble error}) = \sum_{K=\frac{L+1}{2}}^L \binom{L}{K} \epsilon^K (1-\epsilon)^{(L-K)} \quad \left. \vphantom{\sum} \right\} \text{error of bagged ensemble w/ } L \text{ members}$$

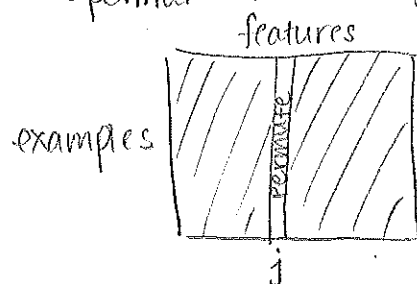
regression:

$$P(\text{ensemble error}) = \left[\frac{1}{L} \sum_{l=1}^L (f(x) + \epsilon_l(x)) \right] - f(x)$$

$\underbrace{\qquad\qquad\qquad}_{\frac{1}{L} \sum_{l=1}^L h_l(x) = h_{\text{bag}}(x)}$

Random Forests

- decision tree of depth D
- \sqrt{d} randomly chosen features; to see which features are most important, do:
 - permute values in j th column



\rightarrow if permuting j has no effect on decision tree...

\rightarrow PRUNE! don't need to split on that feature

Feature Importance cont'd

OOB error of an ensemble: (out of bag)

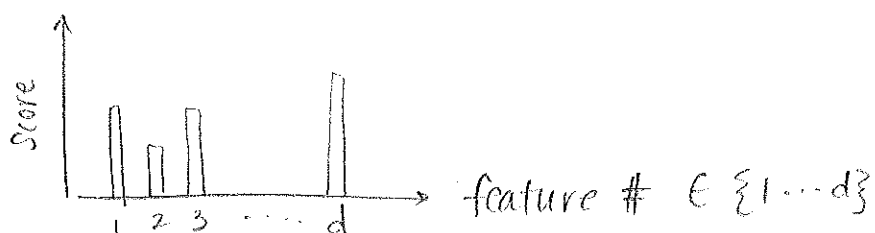
1. cycle over examples
2. find decision trees where example is NOT used (out of bag, OOB)
3. for $i=1$ to m

- use OOB classifier tree to vote on x
- check if majority vote $== y^{(i)}$
- this is OOB error

$$4. \text{Score}(f) = \text{OOBerror}(x) - \text{OOBerror}(x | f \text{ permute})$$

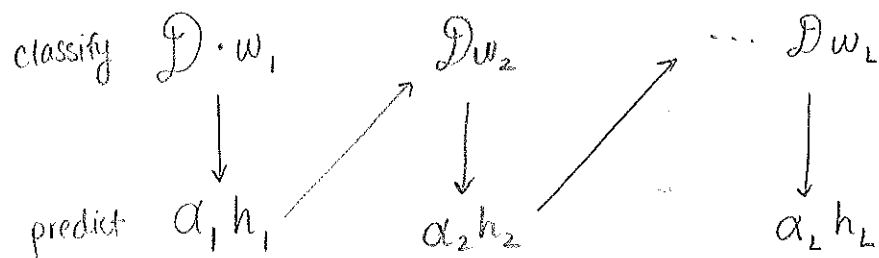
$f \in \{1 \dots d\}$
features

higher score \rightarrow feature more important



Boosting: another method of uncorrelating groups of examples

- associate a weight w / each example



↑ weights on mistakes

↓ weights on correctly classified examples

$$h_{\text{ensemble}}(x) = \text{sgn} \left(\sum_{\ell=1}^L \alpha_{\ell} h_{\ell}(x) \right)$$

weighted vote of h from each classifier

Boosting con't

algorithm:

1. initialize $w_i^{(1)} = \frac{1}{m}$ for $i=1 \dots m$ (all examples have same weight)
2. for $l=1$ to L do:

• learn h_l to minimize $J_l = \frac{1}{m} \sum_{i=1}^m W_l^{(i)} \mathbb{I}(h_l(x^{(i)}) \neq y^{(i)})$

[minimize the # of wrongly classified examples]

• calculate error rate of h_l : $\epsilon_l = \frac{\sum_{i=1}^m W_l^{(i)} \mathbb{I}(h_l(x^{(i)}) \neq y^{(i)})}{\sum_{i=1}^m W_l^{(i)}} \rightarrow \left[\text{this is just normalization} \right]$

* if $\epsilon_l \geq 0.5$, break out of loop! (termination criteria)

• calculate $\alpha_l = \frac{1}{2} \ln \left\{ \frac{1 - \epsilon_l}{\epsilon_l} \right\}$

• update example $1 \dots m$ weights!

$$\rightarrow W_{l+1}^{(i)} = \begin{cases} W_l^{(i)} \exp(-\alpha_l) & \text{if } x^{(i)} \text{ correctly classified} \\ W_l^{(i)} \exp(+\alpha_l) & \text{if } x^{(i)} \text{ incorrectly classified} \end{cases}$$

Working out an example:

10 "examples"

2 iterations

	+	+	-
+		+	
	-	-	
+			-
	-		

$W_l = 0.1$
 $\epsilon_l = 0.3$
 $\alpha_l = \frac{1}{2} \ln \left(\frac{0.7}{0.3} \right) = 0.42$

h_1

updates:

$$W_{l+1} = \begin{cases} W_l e^{-0.42} & \text{if } x \text{ correct} \\ W_l e^{0.42} & \text{if } x \text{ incorrect} \end{cases}$$

	+	+	
			+
+	-	-	
			-
+			
	-		

$\epsilon_l = 0.21$
 $\alpha_l = \frac{1}{2} \ln \left(\frac{0.79}{0.21} \right) = 0.65$

h_2

$$\text{sgn}(\alpha_1 h_1 + \alpha_2 h_2 \dots)$$

Adaboost (notes are spotty b/c I was PMS-ing that day ...)
... aka in a bad mood ☹

$$J = \sum_{i=1}^m \exp(-y^{(i)} H_L(x^{(i)})) \text{ where } H_L(x) = \sum_{j=1}^L \alpha_j h_j(x) \text{ where } h_j(x) \in \{+1, -1\}$$

We keep $\alpha_1, \dots, \alpha_{L-1}, h_1, \dots, h_{L-1}$ constant and learn α_L, h_L to min. J

$$J = \sum_{i=1}^m \underbrace{\exp(-y^{(i)} H_{L-1}(x^{(i)}))}_{w_L^{(i)}} \cdot \exp(-y^{(i)} \alpha_L h_L(x^{(i)}))$$

$$J = \sum_{i=1}^m w_L^{(i)} \exp(-y^{(i)} \alpha_L h_L(x^{(i)}))$$

$$J = \sum_{\substack{i \text{ correctly} \\ \text{classified}}} w_L^{(i)} \exp(-\alpha_L) + \sum_{\substack{i \text{ incorrectly} \\ \text{classified}}} w_L^{(i)} \exp(\alpha_L)$$

|||

$$J = (\exp(+\alpha_L) - \exp(-\alpha_L)) \underbrace{\left[\sum_{i=1}^m w_L^{(i)} \mathbb{I}(h_L(x^{(i)}) \neq y^{(i)}) \right]}_A + \underbrace{\left[\sum_{i=1}^m w_L^{(i)} \right]}_B \exp(-\alpha_L)$$

$$J = (\exp(+\alpha_L) - \exp(-\alpha_L)) A + \exp(-\alpha_L) B \quad \frac{A}{B} \stackrel{\text{def}}{=} \varepsilon_L \text{ what}$$

now take gradient and set to 0!!

$$\frac{\partial J}{\partial \alpha_L} = 0 \implies [\exp(\alpha_L) + \exp(-\alpha_L)] A - B \exp(-\alpha_L) = 0$$

$$[\exp(\alpha_L) + \exp(-\alpha_L)] \varepsilon_L - \exp(-\alpha_L) = 0$$

$$\therefore \varepsilon_L = \frac{\exp(-\alpha_L)}{\exp(\alpha_L) + \exp(-\alpha_L)}$$

$$1 - \varepsilon_L = \frac{\exp(\alpha_L)}{\exp(\alpha_L) + \exp(-\alpha_L)}$$

$$\frac{1 - \varepsilon_L}{\varepsilon_L} = \exp(2\alpha_L)$$

Gradient Boosting (apparently Adaboost is a special case of GB)

regression example:

$$\mathcal{D} = \{(x^{(i)}, y^{(i)}) \mid 1 \leq i \leq m, x^{(i)} \in \mathbb{R}, y^{(i)} \in \mathbb{R}\} \quad f(x) \text{ is true fxn}$$

goal: find $h(x)$ to min $\sum_{i=1}^m L(y^{(i)}, h(x^{(i)}))$

$$\hat{h}(x) = h_1(x) + \underbrace{h_2(x)}_{y - h_1(x) \text{ the residual!!}}$$

$$\rightarrow \mathcal{D}_2 = \{(x^{(i)}, y^{(i)} - h_1(x)) \mid 1 \leq i \leq m, x^{(i)} \in \mathbb{R}^d\}$$

\rightarrow find h_2 to min squared err loss over \mathcal{D}_2

\rightarrow focusing on errors made by last classifier

the gradient is the residual!

$$\text{ex: } J = \frac{1}{2} \sum_{i=1}^m (y^{(i)} - h_1(x^{(i)}))^2 \quad \frac{\partial J}{\partial h_1} = - \sum_{i=1}^m (y^{(i)} - h_1(x^{(i)}))$$

$$h_1 \leftarrow h_1 - \frac{\partial J}{\partial h_1} \quad \text{tada!}$$

Gradient Boost

vs Adaboost

- generates learners during learning process
- 1st learner predicts $y^{(i)}$'s
- 2nd learner predicts loss
- predict loss until threshold is hit!

- users specify a set of learners at start
- it learns weights of how to add learners together to be stronger!
- if learner classifies incorrectly \rightarrow lower its weight

learners = classifiers

$\underbrace{(\text{loss})}_{\text{sum}(ny)}$

Unsupervised Learning

PCA

- project to lower dimensionality
- max variance in data

find $u, u^T u = 1$ that max var in direction

Short proof:

$$\text{var}(y_i) = \frac{1}{m} \sum_{i=1}^m (y_i^{(i)} - \bar{y}_i)^2 \quad \bar{y}_i = u_i^T X$$

$$= \frac{1}{m} \sum_{i=1}^m (u_i^T X^{(i)} - u_i^T \bar{X})^2$$

$$= \frac{1}{m} \sum_{i=1}^m (u_i^T (X^{(i)} - \bar{X}))^2$$

$$= \frac{1}{m} \sum_{i=1}^m (u_i^T (X^{(i)} - \bar{X})(X^{(i)} - \bar{X})^T u_i)$$

$$u_i^T \underbrace{\left(\frac{1}{m} \sum_{i=1}^m (X^{(i)} - \bar{X})(X^{(i)} - \bar{X})^T \right)}_{\text{cov}(X)} u_i$$

$$\underset{u_i}{\text{argmax}} \quad u_i^T S u_i \quad \text{st.} \quad u_i^T u_i = 1$$

$$\mathcal{L} = u_i^T S u_i + \lambda_i (1 - u_i^T u_i)$$

$$\frac{\partial \mathcal{L}}{\partial u_i} = 2 S u_i - 2 \lambda_i u_i = 0$$

$$S u_i = \lambda_i u_i \quad \text{wow! eigenvectors}$$

K-means

- goal: separate into K clusters
- min inter-point distances within clusters

- max distance between cluster means
- min distance between point and mean

$$\arg \min J = \sum_{i=1}^m \sum_{k=1}^K z_k^{(i)} \|x^{(i)} - \mu_k\|^2$$

μ_k cluster mean
 z_k latent var

Sweep across z_k and μ_k 's

Algorithm:

- choose values for μ
- choose clusters \rightarrow assign $x^{(i)}$'s to cluster k (E-step)
 $\min J$ wrt. z w/ μ fixed
- relocate means \rightarrow update value of each μ (M-step)
 $\min J$ wrt. μ w/ z fixed
- repeat til stable

E-step: assign points

- calculate J for each val of k for each point x
- select k w/ smallest J

M-step: relocate means

- each $x^{(i)}$ is independent
- compute means of the data points assigned to k

Details

- slow, $O(mk)$
- local min possible, convg gaurenteed
- run multiple times

- initialize cluster means w/ data points (sample)
- best for convex shapes

- cluster means too close?
 \rightarrow overfitting



Gaussian Mixture Models

- K-means is deterministic

$$z \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

- GMM is probabilistic

$$z \rightarrow \begin{bmatrix} 0.1 & 0.5 & 0.3 & 0.1 \end{bmatrix}$$

$$p(x) = \sum_{k=1}^K \underbrace{P(z=k)}_{\pi_k} \underbrace{P(x|z=k)}_{\text{Gaussian} \sim N(x|\mu_k, \sigma_k)}$$

$$p(x) = \sum_{k=1}^K \pi_k N(x|\mu_k, \sigma_k) \quad \sum_{k=1}^K \pi_k = 1 \quad \text{constraint}$$

K is # Gaussians \rightarrow can sweep over

GMMs are generative models! Given $K, \pi_k, \mu_k, \Sigma_k$

$$p(z=k|x) = \frac{p(x|z=k)p(z=k)}{p(x)}$$

$$\arg \max_k p(z=k|x) = \frac{N(x|\mu_k, \Sigma_k) \pi_k}{\sum_{k=1}^K p(x|z=k)p(z=k)}$$

$$\mathcal{L}(\mathcal{D}; \pi, \mu, \Sigma) = \prod_{i=1}^m P(x^{(i)}; \pi, \mu, \Sigma)$$

$$= \prod_{i=1}^m \sum_{k=1}^K P(x^{(i)} z^{(i)}; \pi, \mu, \Sigma)$$

$$= \prod_{i=1}^m \sum_{k=1}^K \underbrace{P(z^{(i)}=k; \pi)}_{\pi_k} \underbrace{P(x^{(i)}|z^{(i)}=k; \mu, \Sigma)}_{N(x^{(i)}; \mu_k, \Sigma_k)}$$

- Use EM! $\boxed{\pi, \mu, \Sigma}$ assume

\downarrow E-step

$P(z^{(i)}=k|x)$
(pick max label)

\curvearrowright re-estimate π, μ, Σ
M-step

$$\pi_k = \frac{\sum_{i=1}^m \mathbb{I}(z^{(i)}=k)}{m}$$

$$\mu_k = \frac{\sum_{i=1}^m \mathbb{I}(z^{(i)}=k) x^{(i)}}{\sum_{i=1}^m \mathbb{I}(z^{(i)}=k)}$$

$$\sum_k = \frac{\sum_{i=1}^m \mathbb{I}(z^{(i)}=k) (x^{(i)} - \mu_k)^T (x^{(i)} - \mu_k)}{\sum_{i=1}^m \mathbb{I}(z^{(i)}=k)}$$

GMMs. cont'd

"Hard variant" \rightarrow assign z as one k instead of probability

"Soft variant" \rightarrow complete log-likelihood

$$Q(\theta, \theta^{(t-1)}) = E[\ell_c(\mathcal{D}; \theta^{(t-1)})] \quad \text{E-step}$$

\uparrow
old parameters

$$\theta^{(t)} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^{(t-1)}) \quad \text{M-step}$$

1. guess $\theta = [\pi, \mu, \Sigma] \quad k=1 \dots K$

2. $r_k^{(i)} = P(z^{(i)} = k | x^{(i)}; \theta)$

3. $\theta^{(t)} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^{(t-1)})$

for GMMs

$$\begin{aligned} Q(\theta, \theta^{(t-1)}) &= E\left[\sum_{i=1}^m \log P(x^{(i)}, z^{(i)} | \theta^{(t-1)})\right] \\ &= E\left[\sum_{i=1}^m \log \prod_{k=1}^K (\pi_k N(x^{(i)} | \mu_k, \Sigma_k))^{\mathbb{I}(z^{(i)}=k)}\right] \end{aligned}$$

$$\begin{aligned} Q(\theta, \theta^{(t-1)}) &= E\left[\sum_{i=1}^m \sum_{k=1}^K \mathbb{I}(z^{(i)}=k) [\log \pi_k + \log N(x^{(i)} | \mu_k, \Sigma_k)]\right] \\ &= \sum_{i=1}^m \sum_{k=1}^K r_k^{(i)} \log(\pi_k) + \sum_{i=1}^m \sum_{k=1}^K r_k^{(i)} \log(N(x^{(i)} | \mu_k, \Sigma_k)) \end{aligned}$$

$$\frac{\partial Q}{\partial \pi_k} = 0 \quad \frac{\partial Q}{\partial \mu_k} = 0 \quad \frac{\partial Q}{\partial \Sigma_k} = 0$$

you know the drill!