

PERCEPTRON (X, Y, T)

$\theta, \theta_0 = 0$
 for t in T :
 for $i = 1 \dots n$:
 if $y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \leq 0$:
 $\theta = \theta + y^{(i)} x^{(i)}$
 $\theta_0 = \theta_0 + y^{(i)}$

return θ, θ_0
 * IF Prediction doesn't match label, adjust
 Predict = $\text{sign}(\theta \cdot x)$

Arg Percep (X, Y, T):

Same exact if $y^{(i)}(\theta \cdot x^{(i)} + \theta_0) \leq 0$:
 $\theta = \theta + y^{(i)} x^{(i)}$
 $\theta_0 = \theta_0 + y^{(i)}$
 $\bar{\theta} = \bar{\theta} + c y^{(i)} x^{(i)}$
 $\bar{\theta}_0 = \bar{\theta}_0 + c y^{(i)}$
 $c = c - 1/n$ for each T
 return $\bar{\theta}, \bar{\theta}_0$

* Arg all your changed thetas, inc increment ones

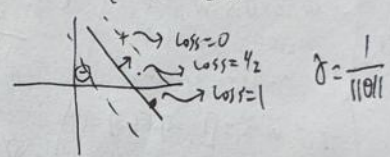
SVM objective: $J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2$

$$\text{Loss}_h(y(\theta \cdot x + \theta_0)) = \max \begin{cases} 0 \\ 1 - y(\theta \cdot x + \theta_0) \end{cases} = \begin{cases} 1 - y(\theta \cdot x + \theta_0) & \text{if } y(\theta \cdot x + \theta_0) \leq 1 \\ 0 & \text{on} \end{cases}$$

Pegasos (X, Y, λ , T)

$\theta = 0$
 for t in T :
 select random i :
 $\eta = 1/t$
 if $y^{(i)} \theta \cdot x^{(i)} \leq 1$:
 $\theta = (1 - \eta \lambda) \theta + \eta y^{(i)} x^{(i)} \leftarrow \theta - \eta \nabla \Rightarrow \nabla = -x y + \lambda \theta$
 else
 $\theta = (1 - \eta \lambda) \theta$

$$\text{Alt H Loss} = \begin{cases} 1 - \frac{\partial(x, y, \theta, \theta_0)}{\partial_{\text{ref}}} & \text{if } \partial(x, y, \theta, \theta_0) < \partial_{\text{ref}} \\ 0 & \text{on} \end{cases}$$



Regression

GD: $x_1 \leftarrow x_0 - \nabla_x f(x_0)$
 NGD: estimate ∇ by $\frac{f(x_0 + \delta^{(i)}) - f(x_0 - \delta^{(i)})}{2\delta}$

$f(x, \theta, \theta_0) = \theta \cdot x + \theta_0 \rightarrow \text{predict}$
 Learning: $\text{Loss} = \frac{(y^{(i)} - \theta \cdot x^{(i)})^2}{2} + \frac{\lambda \|\theta\|^2}{2}$ in $J = \frac{1}{n} \sum L(x, y, \theta, \theta_0) + \frac{\lambda \|\theta\|^2}{2}$

- structural error - when you have linear model for poly solution
- Estimation error = test error. For less dim in feature vectors, ↓
- Regularization prevents overfitting

Set $\theta = 0$
 Random select $i \in [1 \dots n]$
 $\theta^{k+1} = \theta^k + \eta_k (y^{(i)} - \theta^k \cdot x^{(i)}) x^{(i)}$
 Ridge Regression: w/ reg $\Rightarrow \theta^{k+1} = (1 - \lambda \eta_k) \theta^k + \eta_k (y^{(i)} - \theta^k \cdot x^{(i)}) x^{(i)}$
 (w/ regularization)

Matrix Split - Learning

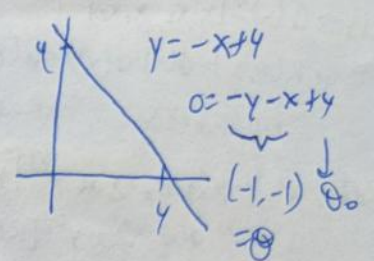
$y_{a,i}$ - given
 B is every v that appears in a_i
 z is all $y_{a,i}$ in a column
 To solve for optimal U 's: set $\nabla J = 0$
 $0 = -(B)^T (z - B U_a) + \lambda \cdot U_a$
 or $0 = -(B)^T (z - B U_a - b_v - b_u^{(a)}) + \lambda U_a$

Ex
 $y_{1,2} = 1$
 $U = [u_1, u_2]$
 solve for $[u_1, u_2]$
 $J(U, V) = \frac{1}{2} \sum [y_{a,i} - [u v]]^2 + \frac{\lambda}{2} (u_1^2 + v_1^2)$
 $0 = \nabla J_U = \frac{1}{2} (2 \cdot (-U) [y_{a,i} - UV]) + \lambda U_1$
 $0 = \nabla J_{u_1} = \lambda U_1 = 0$
 $0 = \nabla J_{u_2} = -u [1 - u v_2] = 0$
 $v_2 = \frac{1}{u}$

Recommendation: Predict $y_{a,i} \rightarrow u_{a,i} v_i = x_{a,i}$

$J = \sum_{i: (a,i) \in D} \frac{(y_{a,i} - u_{a,i} v_i)^2}{2} + \frac{\lambda}{2} \|u\|^2$ for optimizing $u^{(a)}$
 $\frac{\lambda}{2} \|v\|^2$ for optimizing $v^{(i)}$

→ Alternating Optimizing
 w/ w/o offsets: $J(u, v) = \sum \frac{(y_{a,i} - u_{a,i} v_i - b_u^{(a)} - b_v^{(i)})^2}{2} + \frac{\lambda}{2} \sum \|u\|^2 + \frac{\lambda}{2} \sum \|v\|^2$ all b's that showing up in a, i



Geometric margin - if incorrect class, find θ be a c'd
 θ in direc of pos p's

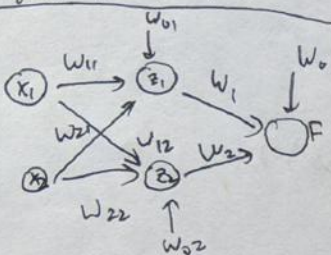
kernel $\phi(x) = [x_1, x_2, x_1 x_2]^T$
 $K(x, x') = \phi(x) \cdot \phi(x') = k(x, x')$
 $K(x, x) \geq 0$ • Radial Basis $\rightarrow e^{-\alpha \|x - x'\|^2/2}$

Predict x' : $\theta \cdot \phi(x) + \theta_0 = \sum a_i y^{(i)} [k(x^{(i)}, x') + 1]$
 or $\theta \cdot \phi(x) = \sum a_i y^{(i)} [k(x^{(i)}, x')]$

Perceptron $\alpha \in [-0.01, 0]$
 IF $y^{(u)} (\sum_{i=1}^n \alpha_i y^{(i)} [k(x^{(i)}, x') + 1]) \leq 0$
 $\alpha \leftarrow \alpha + 1$

Regression: Loss function is to the effect of
 $J(\theta) = \frac{1}{n} \sum_{t=1}^n (y^{(t)} - \theta \cdot \phi(x^{(t)}))^2 + \frac{\lambda}{2} \|\theta\|^2$
 $\vec{\alpha} = \vec{y} (n\lambda I + K)^{-1}$
 $\nabla J(\theta) = -\lambda \sum \phi(x^{(t)}) + \lambda \theta = 0$ Gram Matrix

NN



$$z_i = [w_{i1}, w_{i2}] \cdot [x_1, x_2] + w_{i0}$$

$f(z) = z \rightarrow$ Linear

$f(z) = \text{Relu}(z) \begin{cases} 0 & z \leq 0 \\ z & \text{ow} \end{cases}$

$f(z) = \text{tanh}(z)$

$f(z) = \text{sigmoid}(z) = \frac{1}{1+e^{-z}}$

$f(z) = \text{softmax}(z) = \frac{e^{z_j}}{\sum_k e^{z_k}}$ - outputs a probability dist vector

Learning/Backprop

$a = \text{sm}(z)$

Loss = $(1-y)a$ if $y \leq 1$ or $(1-y)a$ if $y > 1$

usu you take $\frac{\partial \text{Loss}}{\partial \text{weight}}$ or $\frac{\partial \text{Loss}}{\partial w_{02}} = -y \frac{\partial f}{\partial w} = -y \frac{\partial f}{\partial z} \left(\frac{\partial z}{\partial w} \right)$
 $= -y \frac{\partial f}{\partial z} \left(\frac{\partial}{\partial z} \sum w \cdot x \right)$

For simple NN w/ linear act.

$$w_{i1} \leftarrow w_{i1} + \eta y^{(u)} x_1^{(i)}$$

$$w_{i0} \leftarrow w_{i0} + \eta y^{(i)}$$

NLL $\ell(a, y) = -\sum_{j=1}^K y_j \log a_j$ - Objective for probability dist \vec{a}, \vec{y}

$$\frac{\partial}{\partial w_{jk}} \text{NLL}(a^t, y) = x_k(a_j^t - y_j)$$

probability is positive a from Sigmoid function

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

$$y' = y(1-y) = \frac{y}{2z}$$

* larger weights cause sigmoid to saturate

Relu deriv

$$f = \text{Relu}[w \cdot x + w_0]$$

$$\frac{\partial f}{\partial w_i} = x_i [0 \text{ if } z \leq 0, 1 \text{ if } z > 0]$$

CNN - Filters, pooling

where things are what is in each region

CNN w/ 100 units, 10 filters each w/ 5 inputs w/o bias \rightarrow

100 units x 10 feature maps = 1000 output
 \rightarrow needs 5.10+10 weights

Act func

Loss func

hidden units: Relu, tanh, nonlinearity

linear \rightarrow Quadratic
 Sigmoid (prob) \rightarrow NLL
 softmax (prob) \rightarrow NLL
 multiclass

Ex Ridge Regression w/ kernels

$$J(\theta) = \frac{\lambda}{2} \|\theta\|^2 + \frac{1}{n} \sum (y^{(t)} - \theta \cdot \phi(x^{(t)}))^2$$

$$\nabla_{\theta} \rightarrow \lambda \hat{\theta} - \frac{1}{n} \sum (y^{(t)} - \theta \cdot \phi(x^{(t)})) \phi(x^{(t)}) = 0$$

$$J(\theta) \geq J(\hat{\theta}) \text{ for all } \theta \in \mathbb{R}^d$$

* If $\lambda \uparrow, \|\theta\| \downarrow$

* Adding features to $\phi(x)$ will \downarrow training error

Predict: $\hat{\theta} \cdot \phi(x) \begin{cases} 1 & z \geq 0.5 \\ 0 & \text{ow} \end{cases}$

MDP | Horizon shift

$$\pi(s) = \max_{a \in A} [R(s, a) + \sum_{s'} \pi(s, a, s') \cdot V_{t+1}(s')]$$

Val Iteration - Fin

$$Q_0(s, a) = 0$$

For i in $1:H$

For s, a in S, A

$$Q_i(s, a) = R(s, a) + \sum_{s'} \tau(s, a, s') \max_a Q_{i-1}(s', a')$$

$$\text{Inf } Q(s, a) = R(s, a) + \gamma \sum_{s'} \tau(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$$

until $\max_{s, a} \|Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)\| \leq \text{VCS}$

MDP cont

$$Q(s,a) := (1-\alpha)Q(s,a) + \alpha(R + \gamma \max_{a'} Q(s_{t+1}, a'))$$

Value Iteration w/ $R(s,a,s')$

$$Q^*(s,a) = \sum_{s' \in S} T(s,a,s') [R(s,a,s') + \gamma V^*(s')]$$

Goal-reward - init 0, all movements then are random until goal state is found

SSP: init -1, 0 for reward \rightarrow all movements are away from prev bc -1 reward
 \rightarrow can discover faster

Deep Q Learning

• Bellman error: $R(s_t, a_t) + \gamma \max_{a'} [Q(s_{t+1}, a', \theta) - Q(s_t, a_t, \theta)]^2$

• Fitted Q learn:

For i in iters:

for k in trajectories:

$t = \text{episode (policy)}$

experiences, append (random-select t)

for e in experiences:

data, labels = experience

model, fit (data, labels)

Goal labels = $r + \gamma \max_{a'} Q(s', a, \theta)$

$\rightarrow \epsilon$ -greedy

RNN

$$S_t = f_1(W^{sx}x_t + W^{ss}s_{t-1})$$

$$z_t = f_2(W^o s_t)$$

$$\text{let } z_t = \sum_{i=0}^t (0.9^{t-i} + 0.8^{t-i}) x_i$$

Smallest dimensionality: 2

$$W^{ss} = \begin{bmatrix} 0.9 & 0 \\ 0 & 0.8 \end{bmatrix}$$

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

$$x_t = y_{t-1}$$

$$S_t = f_1(W^{sx}x_t + W^{ss}s_{t-1})$$

$$y_t = f_2(W^o s_t)$$

$$W^{ss} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

$$W^{sx} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$W^o = \begin{bmatrix} -1 \\ 2 \\ 3 \end{bmatrix}$$

Ex

$R(s,a) \rightarrow 1$ if $s=1$ 4 states,
 2 if $s=3$ 2 actions a,b
 0 otherwise $\gamma=0.9$

Init $Q = 0 \forall s,a$

Q first \rightarrow Rewards

$[0, 0, 1, 1, 0, 0, 2, 2]$

2nd $Q(0,a) = R(0,a) + \gamma[0.9 \cdot 1 + 0.1 \cdot 0] = 0.81$

Ex

Horizons

$R = -1, 0$ otherwise

Goal: 10, $H=15$, $\gamma=1$ $V^*(1) = -9$

ori: 1 $H=\infty$ $V^*(1) = -9$

If $R=0, 1$ otherwise

$V^*(1) = 0$, $H=1$

$V^*(1) = 6$, $H=15$

$V^*(1) = 0$, $H=\infty$

Ex MDP



$$Q_{S3,1} = 1$$

$$Q_{S4,2} = 1 + 0.5 \cdot 10 = 6$$

$$Q_{S3,3} = 1 + 0.5 \cdot 0.5 \cdot 10 + 0.5 \cdot 10 = 11$$

$$Q_{S3,\infty} = 1 + \sum_{i=1}^{\infty} (0.5^i) \cdot 10 = 11$$

CNNEX

8x8 image - Layer 1

Layer 2: 4 filters w/ 3x3 field, stride 2

Layer 3: Max pooling filter 2x2, stride 2

Layer 4: 4 filters with 3x3x4 field w/ stride 1
 \rightarrow channels

Layer 5: Max Pooling w/ filter size 2x2, stride 2

\rightarrow Weights b/w 1, 2 inc. biases:

$$3 \times 3 \times 4 + 4 = 40$$

\rightarrow Each output is 8x8 for padding

\rightarrow 9 pixels in original image contribute to layer 2 pixel

\rightarrow Output of layer 3 is 4 images 4x4 bc stride 2

\rightarrow 4x4 pixels in original image contribute to one pixel in layer 3

$$\rightarrow 3 \times 3 \times 4 + 4 = 148 \text{ b/w 3 and 4}$$

\rightarrow 2x2 large images from layer 5 bc stride 2

Loss Objective

$$J(\theta, \theta_0) = \underbrace{\frac{1}{n} \sum_{i=1}^n \text{Loss}(x^{(i)}, y^{(i)}, \theta, \theta_0)}_{\text{average loss on training examples}} + \underbrace{\lambda R(\theta, \theta_0)}_{\text{regularizer}}$$

SVM Objective - optimize margin, margin size

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2$$

You can use this w/ perceptron classification

$$\gamma = \frac{y^{(i)} (\theta^T x^{(i)})}{\|\theta\|}$$

margin

Maximize $\gamma = \frac{1}{\|\theta\|}$

↓ minimize

* geometric margin \rightarrow margin on graph from correct points

Hinge Loss

taking the agreement as its argument

$$\text{Loss}_h(y(\theta \cdot x + \theta_0)) = \max\{0, 1 - y(\theta \cdot x + \theta_0)\} = \begin{cases} 1 - y(\theta \cdot x + \theta_0) & \text{if } y(\theta \cdot x + \theta_0) \leq 1 \\ 0, & \text{o.w.} \end{cases} \quad (8)$$

By introducing this loss, we do not exclude any linear classifier from being selected but

\rightarrow takes in same x, y as perceptron, returns optimized decision boundary

Algorithm 1 Pegasos Algorithm (without offset)

```

1: procedure PEGASOS( $\{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}, \lambda, T$ )
2:    $\theta = 0$  (vector)
3:   for  $t = 1, \dots, T$  do
4:     Select  $i \in \{1, \dots, n\}$  at random
5:      $\eta = 1/t$ 
6:     if  $y^{(i)} \theta \cdot x^{(i)} \leq 1$  then
7:        $\theta = (1 - \eta\lambda)\theta + \eta y^{(i)} x^{(i)}$ 
8:     else
9:        $\theta = (1 - \eta\lambda)\theta$    else: just add regularization
10:  return  $\theta$ 

```

$$\theta \leftarrow \theta + \eta \nabla$$

* keep in mind what Loss is for: optimizing θ

- when you evaluate a θ : no regularization

- when you select a θ and optimize, use regularization

- when you select for λ , this is when you're trying to figure out how good your loss func was, and you test it on your test data.