

## 1. SVM

### 1.1. Explain why we can set the margin to c = 1 to derive the SVM formulation?

(a) C is just a constant that scales w (orthogonal vector to decision boundary) and b (offset scalar). Setting c = 1 merely modifies the original mathematical problem into a cleaner form (2c/w l2 norm)

### 1.2. Using Langrangian dual formulation, show that the weight vector can be represented as:

(a) The optimal value of w is a linear combination, projection of the other variables

1. SVM

1.2  $L(w, b, \alpha) = \frac{1}{2} w^T w + \sum_{i=1}^m \alpha_i (1 - y_i (w^T x_i + b))$

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^m \alpha_i y_i x_i = 0$$
$$w = \sum_{i=1}^m \alpha_i y_i x_i$$

### 1.3. Explain why only the data points on the “margin” will contribute to the sum

(a) KKT condition  $\alpha_i g_i(w) = 0$ , not every constraint (data point) is relevant.  
Sparsity structure in solution

### 1.4. Simple SVM by hand

(a) Any value below 1 allows a line to linearly separate the data points  
(b) Yes, the margin from the decision boundary will increase or decrease depending on the location of the support vectors from the boundary

## 2. Optimization

### 2.1. See attachment

2. Optimization

D.1 Show step-by-step mathematical derivation for the gradient of the cost function

$$L(\theta) = \sum_{i=1}^m \left\{ -\log(1 + \exp\{-\theta^T x^i\}) + (y^i - 1)\theta^T x^i \right\}$$

1. split into 2 parts (from lectures)

Part 1

Part 2

2. gradient (from lecture slides)

$$\frac{\partial L}{\partial \theta} = \sum_i (y^i - 1)x^i + \frac{\exp(-\theta^T x^i)x^i}{1 + \exp(-\theta^T x^i)}$$

Part 1 via chain rule

$$\frac{\partial \log(1 + e^{\theta^T x^i})}{\partial \theta} = \frac{\partial \log(1 + e^{\theta^T x^i})}{\partial \theta} \cdot \frac{\partial (1 + e^{\theta^T x^i})}{\partial \theta} = \frac{1}{1 + e^{\theta^T x^i}} \cdot (0 + x^i e^{\theta^T x^i}) = \boxed{\frac{x^i e^{\theta^T x^i}}{1 + e^{\theta^T x^i}}}$$

Part 2

$$\frac{\partial (y^i - 1)\theta^T x^i}{\partial \theta} = (y^i - 1)x^i$$

combine

$$\boxed{\sum_{i=1}^m \frac{x^i e^{\theta^T x^i}}{1 + e^{\theta^T x^i}} + (y^i - 1)x^i}$$

## 2.2. Write pseudo code for gradient descent

From lecture slides:

Update rule

$$x^{t+1} = x^t - \gamma_t \nabla f(x^t)$$

$\gamma_t > 0$  is called the step size or learning rate

1. Set random initialization of parameter theta
2. Set step size  $\gamma_t$
3. Set epsilon tolerance
4. While theta + 1 - theta is greater than epsilon, do theta + step\*(gradient)

Initialize parameter  $\theta^0$

Do

gradient of  $\nabla l(\theta)$   
since we are solving  
maximization

$$\theta^{t+1} \leftarrow \theta^t + \gamma_t \sum_i (y^i - 1) x^i + \frac{\exp(-\theta^t \top x^i)}{1 + \exp(-\theta^t \top x^i)}$$

5. While the  $||\theta^{t+1} - \theta^t|| > \epsilon$

6. Return Theta when complete.

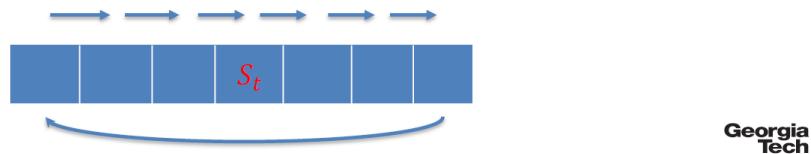
### 2.3. Write pseudo code for Stochastic Gradient Descent

From the lecture slides

- At each iteration, we randomly sample a small subset  $S_t$  data point  $(x_i, y_i), i \in S_t$  (in the extreme case, there is just one sample in  $S_t$ )
- Use gradient estimated using a small subset of data (instead of full data)

$$\nabla \hat{l}(\theta) = \sum_{i \in S_t} (y^i - 1)x^i + \frac{\exp(-\theta^\top x^i)x^i}{1 + \exp(-\theta^\top x^i)}$$

- Each iteration use a different subset  $S_t, t = 1, 2, \dots$
- Eventually loop through the entire training data, and may loop through the data multiple times



3. Initialize random theta
4. Set tolerance + maximum iteration limit
5. While less than maximum iterations + New\_theta is greater than tolerance:
  - 5.1. Randomly sample a subset of data ( $S_t$ )
  - 5.2. New\_Theta = Prev\_Theta + step(1/t)(decreasing) \* gradient(using subset)
6. End if max reached or New\_theta less than tolerance
7. Return New\_Theta

## 2.4 Show Training problem in basic logistic problem is concave

Concave = negative definite

See attachment

l.d  
prove Hessian is negative definite (concave)

$$l(\theta) = \sum \left\{ -\log(1 + \exp\{-\theta^T x^i\}) + (y^i - 1)\theta^T x^i \right\}$$

A) First Derivative :  $\frac{\partial}{\partial \theta} = \sum_{i=1}^m \frac{x^i e^{\theta^T x^i}}{1 + e^{\theta^T x^i}} + (y^i - 1)x^i$   
(from previous problem)

B) Then same 2nd Derivative because Hessian matrix is comprised of 2nd order Derivatives

$$H = \begin{bmatrix} \frac{\partial^2 l}{\partial \theta_1^2} & \cdots & \frac{\partial^2 l}{\partial \theta_m \partial \theta_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 l}{\partial \theta_m \partial \theta_1} & \cdots & \frac{\partial^2 l}{\partial \theta_m^2} \end{bmatrix}$$

2nd Derivative  $\frac{\partial^2 l}{\partial \theta^2} = -\frac{(x^i)^2 e^{\theta^T x^i}}{(e^{\theta^T x^i} + 1)^2}$

Negative definite functions are concave and gradient descent can be used.  
Any local maximum is also a global maximum for concave functions

### 3. Naive Bayes Classifier

#### 3.1

3. Naive Bayes			
1.) class priors		7 total Messages	
spam	not spam	3 spam	
$P(y=0)$	$P(y=1)$	4 non-spam	
3/7	4/7		
Feature Vectors			
$V = \{ \text{secret, offer, low, price, ..., pizza} \}$			
$\begin{matrix} n_1 & n_2 & n_3 & n_4 & n_5 \\ 1 & 2 & 3 & 4 & 5 \end{matrix}$			6 7 8 9 10 11 12 13 14 15
SPAM		Million Dollar offer	
Secret offer today		0 1 0 0 0 0 0 1 1 0 0 0 0 0 0 0	
Secret is secret		1 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0	
Low price for valued customers		2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	
Play secret sports today		0 0 1 1 1 1 0 0 0 0 0 1 0 0 0 0	
Sports is healthy		1 0 0 0 0 0 0 1 0 0 0 1 0 1 0 0	
Low price pizza		0 0 0 0 0 0 0 0 0 0 0 1 1 0 1 0	
		0 0 1 1 0 0 0 0 0 0 0 0 0 0 0 1	

### 3.2

3. Naive Bayes  
 2. Source: Alfie Chen weekly OH

$0 \leq \theta_{c,k} \leq 1 = \text{probability of } k \text{ appear in } c$   
 $C = \text{class } \{0, 1\}$   
 $k = \text{word } (i=1 \text{ to } 15)$

$$\mathcal{L}(\theta, 0) = l(\theta) + \lambda g(\theta)$$

$$\mathcal{L}(\theta, 0) = \left( \sum_{i=1}^m \sum_{k=1}^d x_k^i \log \theta_{y_i, k} \right) \left[ + \lambda \left( \sum_{k=1}^d \theta_{0,k} - 1 \right) \right] + \lambda \left( \sum_{k=1}^d \theta_{1,k} - 1 \right) \lambda g(\theta)$$

first order derivative:

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta_{c,k}} = 0 = \frac{1}{\theta_{c,k}} \sum_{i=1}^m x_k^i + \lambda_{c_{(0,1)}}$$

Then solve for Theta  $\theta_{c,k}$

$$\theta_{c,k} = -\frac{\sum_{i=1}^m x_k^i}{\lambda_{c_{(0,1)}}} = \frac{\# \text{of } k \text{ words}}{\# \text{total words in class } c}$$

a)  $\theta_{0,1} = \frac{3}{9} = \frac{1}{3}$       b)  $\theta_{0,7} = \frac{1}{9}$

c)  $\theta_{1,1} = \frac{1}{15}$       d)  $\theta_{1,5} = \frac{1}{15}$

### 3.3

Naive Bayes

3.3 "today is Secret"

Sources: Lecture Slides

Prior Probability:  $p(y=0) = \frac{4}{7}$     $p(y=1) = \frac{3}{7}$

Posterior  $p(y=c_0|x) = \frac{p(x|y=0)p(y=0)}{p(x|y=0)p(y=0) + p(x|y=1)p(y=1)}$

Today =  $v_7$     $\theta_{0,7} = \frac{1}{9}$     $\theta_{1,7} = \frac{1}{15}$   
is =  $v_{11}$     $\theta_{0,11} = \frac{1}{9}$     $\theta_{1,11} = \frac{1}{15}$   
Secret =  $v_1$ ,    $\theta_{0,1} = \frac{3}{7} = \frac{1}{3}$     $\theta_{1,1} = \frac{1}{15}$

Posterior (via Python)  $\approx >.94$ , so "today is secret" is classified as spam with high probability, assuming .5 cutoff