# Online Learning: LMS and Perceptrons

Partially adapted from slides by Ryan Gabbard and Mitch Marcus (and lots original slides by Lyle Ungar)

Note: supplemental material for today is supplemental; not required!

## Why do online learning?

#### • Batch learning can be expensive for big datasets

• How expensive is it to compute (X<sup>T</sup>X)<sup>-1</sup> for X?



# Why do online learning?

#### Batch learning can be expensive for big datasets

- How hard is it to compute  $(X^TX)^{-1}$ ?
  - np<sup>2</sup> to form X<sup>T</sup>X
  - p<sup>3</sup> to invert
- Tricky to parallelize inversion

#### Online methods are easy in a map-reduce environment

 They are often clever versions of stochastic gradient descent Have you seen map-reduce/hadoop?
A) Yes B) No

## **Online linear regression**

- ♦ Minimize Err = ∑<sub>i</sub> (y<sub>i</sub> w<sup>T</sup>x<sub>i</sub>)<sup>2</sup> using stochastic gradient descent
  - Look at each observation (x<sub>i</sub>,y<sub>i</sub>) sequentially and decrease its error Err<sub>i</sub> = (y<sub>i</sub> - w<sup>T</sup>x<sub>i</sub>)<sup>2</sup>
- LMS (Least Mean Squares) algorithm
  - $\mathbf{w}_{i+1} = \mathbf{w}_i \eta/2 \ dErr_i/d\mathbf{w}_i$
  - $dErr_i/dw_i = -2 (y_i w_i^T x_i) x_i = -2 r_i x_i$  $w_{i+1} = w_i + \eta r_i x_i$

 $r_i X_i$  How do you pick the "learning rate"  $\eta$ ?

Note that *i* is the index for both the iteration and the observation, since there is one update per observation

## **Online linear regression**

LMS (Least Mean Squares) algorithm

 $\mathbf{w}_{i+1} = \mathbf{w}_i + \eta \mathbf{r}_i \mathbf{x}_i$ 

### $\blacklozenge$ Converges for $~~0 < \eta < \lambda_{max}$

- Where  $\lambda_{max}$  is the largest eigenvalue of the covariance matrix  $\bm{X}^T\bm{X}$ 

• Convergence rate is inversely proportional to  $\lambda_{max}/\lambda_{min}$  (ratio of extreme eigenvalues of X<sup>T</sup>X)

## **Online learning methods**

### Least mean squares (LMS)

• Online regression -- L<sub>2</sub> error

#### Perceptron

• Online SVM -- Hinge loss

# **Perceptron Learning Algorithm**

Input: A list T of training examples  $\langle \vec{x}_0, y_0 \rangle \dots \langle \vec{x}_n, y_n \rangle$  where  $\forall i : y_i \in \{+1, -1\}$ Output: A classifying hyperplane  $\vec{w}$ Randomly initialize  $\vec{w}$ ; while model  $\vec{w}$  makes errors on the training data do for  $\langle \vec{x}_i, y_i \rangle$  in T do Let  $\hat{y} = sign(\vec{w} \cdot \vec{x}_i)$ ; if  $\hat{y} \neq y_i$  then  $\vec{w} = \vec{w} + y_i \vec{x}_i$ ; end end end What do we do if error is zero?

Of course, this only converges for linearly separable data

## **Perceptron Learning Algorithm**

For each observation  $(y_i, \mathbf{x}_i)$ 

 $\mathbf{w}_{i+1} = \mathbf{w}_i + \eta \mathbf{r}_i \mathbf{x}_i$ 

Where  $r_i = y_i - sign(\mathbf{w}_i^T \mathbf{x}_i)$ and  $\eta = \frac{1}{2}$ I.e., if we get it right: *no change* if we got it wrong:  $\mathbf{w}_{i+1} = \mathbf{w}_i + y_i \mathbf{x}_i$ 

### **Perceptron Update**



### **Perceptron Update Example II**



# **Properties of the Simple Perceptron**

#### You can prove that

- If it's possible to separate the data with a hyperplane (i.e. if it's linearly separable), then the algorithm will converge to that hyperplane.
- And it will converge such that the number of mistakes M it makes is bounded by

 $M < R^2/\gamma^2$ 

where (assume the true **w** has been normalized:  $||\mathbf{w}^*||_2=1$ )

 $R = \max_i ||\mathbf{x}_{ij}|_2$ size of biggest  $\mathbf{x}$  $\gamma <= y_i \ \mathbf{w}^{*T} \mathbf{x}_i$ > 0 if separable

## **Properties of the Simple Perceptron**

#### But what if it isn't separable?

• Then perceptron is unstable and bounces around

## **Voted Perceptron**

- Works just like a regular perceptron, except you keep track of all the intermediate models you created
- When you want to classify something, you let each of the many models vote on the answer and take the majority

Often implemented after a "burn-in" period

## **Properties of Voted Perceptron**

### ♦ Simple!

#### Much better generalization performance than regular perceptron

- Almost as good as SVMs
- Can use the 'kernel trick'

#### Training is as fast as regular perceptron

- But run-time is slower
  - Since we need **n** models

### **Averaged Perceptron**

- Return as your final model the average of all your intermediate models
- Approximation to voted perceptron
- Again extremely simple!
  - And can use kernels

 Nearly as fast to train and exactly as fast to run as regular perceptron

## Many possible Perceptrons

- If point x<sub>i</sub> is misclassified
  - $\mathbf{w}_{i+1} = \mathbf{w}_i + \eta y_i \mathbf{x}_i$
- Different ways of picking learning rate η
- Standard perceptron:  $\eta = 1$ 
  - Guaranteed to converge to the correct answer in a finite time if the points are separable (but oscillates otherwise)
  - Can get bounds on error even for non-separable case
- Alternate: pick η to maximize the margin (w<sub>i</sub><sup>T</sup>x<sub>i</sub>) in some fashion

## Can we do a better job of picking $\eta$ ?

### Perceptron:

For each observation  $(y_i, \mathbf{x}_i)$   $\mathbf{w}_{i+1} = \mathbf{w}_i + \eta r_i \mathbf{x}_i$ where  $r_i = y_i - sign(\mathbf{w}_i^T \mathbf{x}_i)$ and  $\eta = \frac{1}{2}$ 

Let's use the fact that we are actually trying to minimize a loss function

# **Passive Aggressive Perceptron**

- Minimize the hinge loss at each observation
  - $L(\mathbf{w}_i; \mathbf{x}_i, y_i) = 0$  if  $y_i \mathbf{w}_i^T \mathbf{x}_i \ge 1$  (loss 0 if correct with margin > 1)  $1 - y_i \mathbf{w}_i^T \mathbf{x}_i$  else
- Pick w<sub>i+1</sub> to be as close as possible to w<sub>i</sub> while still setting the hinge loss to zero
  - If point x<sub>i</sub> is correctly classified with a margin of at least 1
    - no change
  - Otherwise
    - $\mathbf{w}_{i+1} = \mathbf{w}_i + \eta y_i \mathbf{x}_i$
    - where  $\eta = L(\mathbf{w}_i; \mathbf{x}_i, y_i) / ||\mathbf{x}_i||^2$
- Can prove bounds on the total hinge loss

### **Passive-Aggressive = MIRA**

$$w_{i+1} = w_i + \frac{y_i - w_i \cdot x_i}{\|x_i\|^2} x_i$$



# Margin-Infused Relaxed Algorithm (MIRA)

#### Multiclass; each class has a prototype vector

- Note that the prototype *w* is like a feature vector *x*
- Classify an instance by choosing the class whose prototype vector is most similar to the instance
  - Has the greatest dot product with the instance

 During training, make the 'smallest' change to the prototype vectors which guarantees correct classification by a specified margin

• "passive aggressive"

# What you should know

### LMS

• Online regression

### Perceptrons

- Online SVM
  - Large margin / hinge loss
- Has nice mistake bounds (for separable case): see wiki
- In practice, use averaged perceptrons
- Passive Aggressive perceptrons and MIRA
  - Change *w* just enough to set its hinge loss to zero.

What we didn't cover: feature selection