

# Stochastic Gradient Descent

Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 14, 2014

©Carlos Guestrin 2005-2014

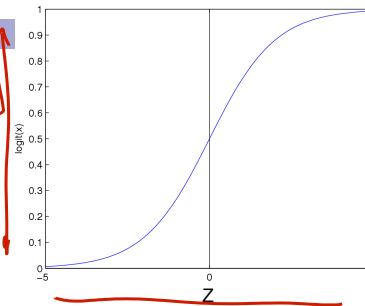
1

## Logistic Regression

Logistic function (or Sigmoid):  $\frac{1}{1 + \exp(-z)}$

- Learn  $P(Y|\mathbf{X})$  directly
  - Assume a particular functional form for link function
  - Sigmoid applied to a linear function of the input features:  $[0,1]$

$$P(Y = 0 | X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$



$$P(Y = 1 | X, W) = 1 - P(Y = 0 | X, W) = \frac{e^{w_0 + \sum_i w_i X_i}}{1 + e^{w_0 + \sum_i w_i X_i}} \leftarrow h_i(x)$$

$X_i$

indicators:  $X_i = \mathbb{1}(\text{county} = \text{USA})$

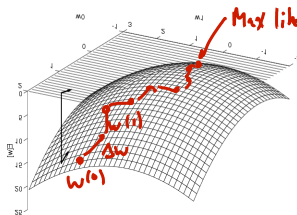
**Features can be discrete or continuous!**

©Carlos Guestrin 2005-2014

2

# Optimizing concave function – Gradient ascent

- Conditional likelihood for Logistic Regression is concave. Find optimum with gradient ascent



Gradient:  $\nabla_w l(\mathbf{w}) = \left[ \frac{\partial l(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial l(\mathbf{w})}{\partial w_n} \right]^T$

Step size,  $\eta > 0$

Update rule:  $\Delta \mathbf{w} = \eta \nabla_w l(\mathbf{w})$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

(choice of  $\eta$ ? from theory)  
 $\eta = \frac{\alpha}{\sqrt{t}}$   
 or  
 $\eta = \frac{\alpha}{t}$   
 or  
 $\eta = \alpha$

- Gradient ascent is simplest of optimization approaches
  - e.g., Conjugate gradient ascent can be much better  
 Newton, LBFGS, ...

for HU

# Gradient Ascent for LR

initialize  $w^{(0)} = 0$  or something else  $\rightarrow$  all converge to same place, concavity

Gradient ascent algorithm: iterate until change  $< \epsilon$

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta \sum_{j=1}^N [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})]$$

small step  $\downarrow$  gradient direction

For  $i=1, \dots, k$ ,

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_{j=1}^N x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})]$$

repeat

# The Cost, The Cost!!! Think about the cost...

*k features  
N data points*

- What's the cost of a gradient update step for LR???

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \sum_{j=1}^N x_i^j [y^j - \hat{P}(Y=1 | \mathbf{x}^j, \mathbf{w}^{(t)})] \right\}$$

*Same for all i → cache*

*for  $i = 0 \dots k$   
total is  $O(Nk^2)$*

*$O(k)$*

*$O(Nk)$*

*With caching  $\hat{P}$ , cost is  $O(Nk)$   $N$  is really large*

*gradient update is  
Vlry expensive for only  
an  $\eta$  step*

# Learning Problems as Expectations

- Minimizing loss in training data:

- Given dataset:  $X^1, X^2, \dots, X^N$ 
  - Sampled iid from some distribution  $p(\mathbf{x})$  on features:
- Loss function, e.g., hinge loss, logistic loss,...
- We often minimize loss in training data:

*$x^j \sim \text{iid } p(\mathbf{x})$*

$$\ell_{\mathcal{D}}(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N \ell(\mathbf{w}, \mathbf{x}^j)$$

*LR.  $-\ln P(y^j | \mathbf{x}^j, \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$*

- However, we should really minimize expected loss on all data:

$$\ell(\mathbf{w}) = E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = \int p(\mathbf{x}) \ell(\mathbf{w}, \mathbf{x}) d\mathbf{x}$$

*expected loss over infinite  $N$*

- So, we are approximating the integral by the average on the training data

discrete.  $\nabla E[f] = \nabla \sum_x f(x) p(x) = \sum_x \nabla f(x) p(x) = E[\nabla f(x)]$

## Gradient descent in Terms of Expectations

- “True” objective function:

$$\ell(\mathbf{w}) = E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = \int p(\mathbf{x}) \ell(\mathbf{w}, \mathbf{x}) d\mathbf{x}$$

- Taking the gradient:

$$\nabla_{\mathbf{w}} \ell(\mathbf{w}) = \nabla_{\mathbf{w}} (E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})]) = E_{\mathbf{x}} [\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x})]$$

- “True” gradient descent rule:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta E_{\mathbf{x}} [\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x})]$$

- How do we estimate expected gradient?

©Carlos Guestrin 2005-2014

7

## SGD: Stochastic Gradient Ascent (or Descent)

- “True” gradient:  $\nabla \ell(\mathbf{w}) = E_{\mathbf{x}} [\nabla \ell(\mathbf{w}, \mathbf{x})]$

- Sample based approximation:

$$\nabla \ell(\mathbf{w}) = E_{\mathbf{x}} [\nabla \ell(\mathbf{w}, \mathbf{x})] \approx \frac{1}{N} \sum_{j=1}^N \nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x}^j)$$

approx. true gradient by sample gradient over dataset

- What if we estimate gradient with just one sample???

- Unbiased estimate of gradient  $E_{\mathbf{x}} [\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x}^t)] = \nabla \ell(\mathbf{w})$
- Very noisy!
- Called stochastic gradient ascent (or descent)
  - Among many other names
- VERY useful in practice!!!

©Carlos Guestrin 2005-2014

8

# Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

$$E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = E_{\mathbf{x}} \left[ \ln P(y|\mathbf{x}, \mathbf{w}) - \lambda \|\mathbf{w}\|_2^2 \right]$$

- Batch gradient ascent updates: *batch approach uses all data*

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \frac{1}{N} \sum_{j=1}^N x_i^{(j)} [y^{(j)} - P(Y=1|\mathbf{x}^{(j)}, \mathbf{w}^{(t)})] \right\}$$

*O(Nk) per iteration*

- Stochastic gradient ascent updates: *pick a next data point  $x^{(t)}$ ,  $y^{(t)}$*

- Online setting: *Often, in practice, mini-batches: pick, typically 10-100 data points*

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y=1|\mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

*cost per iteration*  
*O(k)*

*do an update as if  $x^t, y^t$  were all the data*

©Carlos Guestrin 2005-2014

9

# Stochastic Gradient Descent: general case

- Given a stochastic function of parameters:  $f(w) = E_{\mathbf{x}} [f(w, \mathbf{x})]$

- Want to find maximum

$$w^* \leftarrow \underset{w}{\operatorname{argmin}} f(w) = \underset{w}{\operatorname{argmin}} E_{\mathbf{x}} [f(w, \mathbf{x})]$$

- Start from  $\mathbf{w}^{(0)}$  e.g.  $w^{(0)} = 0$

- Repeat until convergence:

- Get a sample data point  $\mathbf{x}^t$
- Update parameters:

$$w^{(t+1)} \leftarrow w^{(t)} - \eta_t \nabla_w f(w, \mathbf{x}^t)$$

- Works on the online learning setting!
- Complexity of each gradient step is constant in number of examples!
- In general, step size changes with iterations

*$\eta_t$  decreases with iterations, from theory, typically:  $\eta_t = \frac{\alpha}{t}$   
 $\eta_t = \frac{\alpha}{\sqrt{t}}$*

©Carlos Guestrin 2005-2014

10

## What you should know...

- Classification: predict discrete classes rather than real values
- Logistic regression model: Linear model
  - Logistic function maps real values to  $[0, 1]$
- Optimize conditional likelihood
- Gradient computation
- Overfitting
- Regularization
- Regularized optimization
- Cost of gradient step is high, use stochastic gradient descent

©Carlos Guestrin 2005-2014

11

# Boosting

Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 14, 2014

©Carlos Guestrin 2005-2014

12

## Fighting the bias-variance tradeoff

- **Simple (a.k.a. weak) learners are good**
  - e.g., naïve Bayes, logistic regression, decision stumps (or shallow decision trees)
  - Low variance, don't usually overfit too badly
- **Simple (a.k.a. weak) learners are bad**
  - High bias, can't solve hard learning problems
- Can we make weak learners always good???
- **No!!!**
- **But often yes...**

©Carlos Guestrin 2005-2014

13

## Voting (Ensemble Methods)

- Instead of learning a single (weak) classifier, learn **many weak classifiers** that are **good at different parts of the input space**
  - **Output class:** (Weighted) vote of each classifier
    - Classifiers that are most "sure" will vote with more conviction
    - Classifiers will be most "sure" about a particular part of the space
    - On average, do better than single classifier!
- $h_i: X \rightarrow Y \in \{-1, +1\}$
- $t^{\text{th}}$  classifier
- $$H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$$
- vote weight
- e.g.,  $h_t(x) = \begin{cases} +1 & \text{if } x_i = 1 \\ -1 & \text{if } x_i = 0 \end{cases}$  if email has word 'CE546'  $\Rightarrow$  not spam (-1)  
else  $\Rightarrow$  spam (+1)
- **But how do you ???**
    - force classifiers to learn about different parts of the input space?
    - weigh the votes of different classifiers?

©Carlos Guestrin 2005-2014

14

# Boosting [Schapire, 1989]

- Idea: given a weak learning alg, run it multiple times on (reweighted) training data, then let learned classifiers vote
- On each iteration  $t$ :
  - weight each training example by how incorrectly it was classified *thus far*
  - Learn a hypothesis –  $h_t$
  - A strength for this hypothesis –  $\alpha_t$
- Final classifier:  $H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$
- **Practically useful**
- **Theoretically interesting**

©Carlos Guestrin 2005-2014

15

# Learning from weighted data

- **Sometimes not all data points are equal**
  - Some data points are more equal than others
- **Consider a weighted dataset**
  - $D(j)$  – weight of  $j$ th training example  $(x^j, y^j)$
  - Interpretations:
    - $j$ th training example counts as  $D(j)$  examples
    - If I were to “resample” data, I would get more samples of “heavier” data points
- **Now, in all calculations, whenever used,  $j$ th training example counts as  $D(j)$  “examples”**

for gradient descent

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \sum_{j=1}^N D(j) \nabla_w \ell(w, x^j)$$

data point is weighted

©Carlos Guestrin 2005-2014

16



# AdaBoost

Goal is to focus  $h_t$  on hard points, ones that are not classified well, by increasing their weight.

- Initialize weights to uniform dist:  $D_1(j) = 1/N$
- For  $t = 1 \dots T$ 
  - Train weak learner  $h_t$  on distribution  $D_t$  over the data
  - Choose weight  $\alpha_t$
  - Update weights: for each point  $j$

$$D_{t+1}(j) = \frac{D_t(j) \exp(-\alpha_t y^j h_t(x^j))}{Z_t}$$

$\alpha_t > 0$  magic, next slide  
for each point  $j$

Where  $Z_t$  is normalizer:  
re normalize so  $D_t(j)$  always adds up to 1

$$Z_t = \sum_{j=1}^N D_t(j) \exp(-\alpha_t y^j h_t(x^j))$$

uniform  
focus on parts with high weight  
old weight  
new weight  
 $y^j h_t(x^j) > 0 \Rightarrow$  correct class  
 $\Rightarrow \exp(-\alpha_t y^j h_t(x^j)) < 1 \Rightarrow$  weight decrease  
 $y^j h_t(x^j) < 0 \Rightarrow$  incorrect class  
 $\Rightarrow \exp(-\alpha_t y^j h_t(x^j)) > 1 \Rightarrow$  weight increase

Output final classifier:

$$H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$$

## Picking Weight of Weak Learner

- Weigh  $h_t$  higher if it did well on training data (weighted by  $D_t$ ):

Magic:  $\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$

if  $\epsilon_t = 0 \Rightarrow h_t$  perfect on weighted data  
 $\Rightarrow h_t$  also perfect on original data  
 $\Rightarrow \alpha_t = \infty$

$\frac{1}{2} < \epsilon_t < 1$ ,  $h_t$  worse than random, but  $-h_t$  is better than random  
 $\Rightarrow \alpha_t < 0 \Rightarrow$  flip  $h_t$

$\epsilon_t = \frac{1}{2} \Rightarrow h_t$  is same as random guess  
 $\Rightarrow \alpha_t = 0 \Rightarrow$  ignore  $h_t$

- Where  $\epsilon_t$  is the weighted training error:

$$\epsilon_t = \sum_{j=1}^N D_t(j) \mathbb{1}[h_t(x^j) \neq y^j]$$

# Why choose $\alpha_t$ for hypothesis $h_t$ this way?

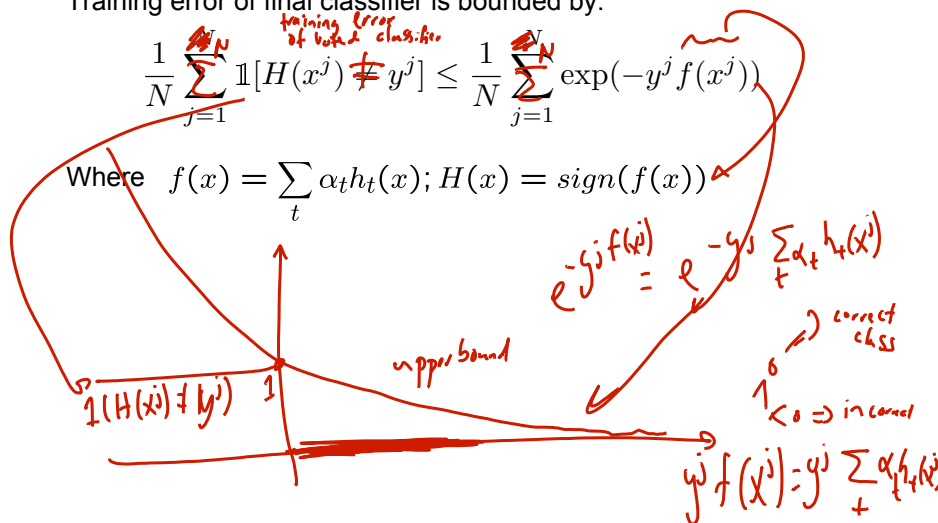
[Schapire, 1989]

Training error of final classifier is bounded by:

$$\frac{1}{N} \sum_{j=1}^N \mathbb{1}[H(x^j) \neq y^j] \leq \frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j))$$

*training error of both classifiers*

Where  $f(x) = \sum_t \alpha_t h_t(x)$ ;  $H(x) = \text{sign}(f(x))$



©Carlos Guestrin 2005-2014

19

# Why choose $\alpha_t$ for hypothesis $h_t$ this way?

[Schapire, 1989]

Training error of final classifier is bounded by:

$$\frac{1}{N} \sum_{j=1}^N \mathbb{1}[H(x^j) \neq y^j] \leq \frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j)) = \prod_{t=1}^T Z_t$$

*training error*      *exponential upper bound*

Where  $f(x) = \sum_t \alpha_t h_t(x)$ ;  $H(x) = \text{sign}(f(x))$



Why? see homework magic of telescopic if  $Z_t < 1$

$\Rightarrow$  error bound decreases exponentially with  $t$

©Carlos Guestrin 2005-2014

20

## Why choose $\alpha_t$ for hypothesis $h_t$ this way?

[Schapire, 1989]

Training error of final classifier is bounded by:

$$\frac{1}{N} \sum_{j=1}^N \mathbb{1}[H(x^j) \neq y^j] \leq \frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j)) = \prod_{t=1}^T Z_t$$

Where  $f(x) = \sum_t \alpha_t h_t(x)$ ;  $H(x) = \text{sign}(f(x))$

**If we minimize  $\prod_t Z_t$ , we minimize our training error**

AdaBoost tightens this bound greedily, by choosing  $\alpha_t$  and  $h_t$  on each iteration to minimize  $Z_t$ .

$$Z_t = \sum_{j=1}^N D_t(j) \exp(-\alpha_t y^j h_t(x^j))$$