

Case Study 1: Estimating Click Probabilities

Stochastic Gradient Descent (continued)

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Carlos Guestrin
January 17th, 2013

©Carlos Guestrin 2013

1

What is the Perceptron Doing???

- When we discussed logistic regression:
 - Started from maximizing conditional log-likelihood
- When we discussed the Perceptron:
 - Started from description of an algorithm
- What is the Perceptron optimizing????

©Carlos Guestrin 2013

2

Perceptron & Stochastic Gradient descent

- Perceptron update:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right] y^{(t)} \mathbf{x}^{(t)}$$

Complexity in terms of N: $O(1)$

$O(N)$

1 sample estimate of $E[\nabla \ell(\mathbf{w})]$
Stochastic Gradient descent!!
with $\eta=1$

- Batch hinge minimization update:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \frac{1}{N} \sum_{i=1}^N \left\{ \mathbb{1} \left[y^{(i)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(i)}) \leq 0 \right] y^{(i)} \mathbf{x}^{(i)} \right\}$$

N sample estimate of $E[\nabla \ell(\mathbf{w})]$

©Carlos Guestrin 2013

3

Stochastic Gradient Descent: general case

after sufficiently large t , running avg of loss approx. true loss, use for stopping -

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

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

- Start from $\mathbf{w}^{(0)}$; e.g., $\mathbf{w}^{(0)} = \mathbf{0}$
- Repeat until convergence:

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

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

Stopping criteria
very hard...
→ theory bounds # iterations in terms of uncomputable constants

- Works on the online learning setting!
- Complexity of gradient computation is constant in number of examples!

→ in practice, validation set

- In general, step size changes with iterations, e.g., $\eta_t = \frac{\kappa}{t}$ for $\kappa > 0$

©Carlos Guestrin 2013

4

Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

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

- Batch gradient ascent updates:

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

- Stochastic gradient ascent updates:

- Online setting:

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

↑ 1 point at a time

average over N datapoints

©Carlos Guestrin 2013

5

Convergence rate of SGD

- Theorem:

- (see Nemirovski et al '09 from readings)
- Let f be a strongly convex stochastic function
- Assume gradient of f is Lipschitz continuous and bounded

- Then, for step sizes:

- The expected loss decreases as $O(1/t)$:

©Carlos Guestrin 2013

6

Convergence rates for gradient descent/ascent versus SGD

- Number of Iterations to get to accuracy

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \leq \epsilon$$

- Gradient descent:
 - If func is strongly convex: $O(\ln(1/\epsilon))$ iterations
- Stochastic gradient descent:
 - If func is strongly convex: $O(1/\epsilon)$ iterations
- Seems exponentially worse, but much more subtle:
 - Total running time, e.g., for logistic regression:
 - Gradient descent:
 - SGD:
 - SGD can win when we have a lot of data
 - And, when analyzing true error, situation even more subtle... expected running time about the same, see readings

©Carlos Guestrin 2013

7

What you need to know

- Perceptron is optimizing hinge loss
- Subgradients and hinge loss
- (Sub)gradient descent for hinge objective
- Objective functions in ML as expectations
- Gradient estimation, rather than objective estimation
- Stochastic gradient descent -> estimate gradient from single training example
 - Mini-batches possible and useful
- Stochastic gradient ascent for logistic regression
- Analysis of stochastic gradient descent
 - Decreasing step size fundamental here
- Comparing analysis of stochastic gradient descent with gradient descent

©Carlos Guestrin 2013

8

Case Study 1: Estimating Click Probabilities

Tackling an Unknown Number of Features with Sketching

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Carlos Guestrin
January 17th, 2013

©Carlos Guestrin 2013

9

Problem 1: Complexity of Update Rules for LR and Perceptron

- Perceptron update:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right] y^{(t)} \mathbf{x}^{(t)}$$

- Logistic regression update:

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

- Complexity of updates:

- Constant in number of data points
- In number of features?
 - Problem both in terms of computational complexity and sample complexity

- What can we with very high dimensional feature spaces?

- Kernels not always appropriate, or scalable
- What else?

©Carlos Guestrin 2013

10

Problem 2: Unknown Number of Features

- For example, bag-of-words features for text data:
 - “Mary had a little lamb, little lamb...”

- What’s the dimensionality of \mathbf{x} ?
- What if we see new word that was not in our vocabulary?
 - Obamacare

 - Theoretically, just keep going in your learning, and initialize $\mathbf{w}_{\text{Obamacare}} = 0$
 - In practice, need to re-allocate memory, fix indices,... A big problem for Big Data

©Carlos Guestrin 2013

11

What Next?

- Hashing & Sketching!
 - Addresses both dimensionality issues and new features in one approach!

- Let’s start with a much simpler problem: Is a string in our vocabulary?
 - Membership query
- How do we keep track?
 - Explicit list of strings
 - Very slow

 - Fancy Trees and Tries
 - Hard to implement and maintain

 - Hash tables?

©Carlos Guestrin 2013

12

Hash Functions and Hash Tables

- Hash functions map **keys** to integers (bins):
 - Keys can be integers, strings, objects,...

- Simple example: **mod**
 - $h(i) = (a \cdot i + b) \% m$

 - Random choice of (a,b) (usually primes)
 - If inputs are uniform, bins are uniformly used
 - From two results can recover (a,b), so not pairwise independent -> Typically use fancier hash functions
- Hash table:
 - Store list of objects in each bin
 - Exact, but storage still linear in size of object ids, which can be very long
 - E.g., hashing very long strings, entire documents

©Carlos Guestrin 2013

13

Hash Bit-Vector Table-based Membership Query

- Approximate queries with one-sided error: Accept false positives only
 - If we say no, element is not in set
 - If we say yes, element is very to be likely in set

- Given hash function, keep binary bit vector **v** of length *m*:

- Query Q(i): Element *i* in set?
 -
 -

- Collisions:

- Guarantee: One-sided errors, but may make many mistakes
 - How can we improve probability of correct answer?

©Carlos Guestrin 2013

14

Bloom Filter: Multiple Hash Tables

- Single hash table -> Many false positives
- Multiple hash tables with independent hash functions
 - Apply $h_1(i), \dots, h_d(i)$, set all bits to 1
- Query $Q(i)$?
- Significantly decrease probability of false positives

©Carlos Guestrin 2013

15

Analysis of Bloom Filter

- Want to keep track of n elements with false positive probability of $\delta > 0$... how large m & d ?
- Simple analysis yields:

$$m = \frac{n \log_2 \frac{1}{\delta}}{\ln 2} \approx 1.5n \log_2 \frac{1}{\delta}$$

$$d = \log_2 \frac{1}{\delta}$$

©Carlos Guestrin 2013

16

Sketching Counts

- Bloom Filter is super cool, but not what we need...
 - We don't just care about whether a feature existed before, but to keep track of counts of occurrences of features!

- Recall Perceptron update:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right] y^{(t)} \mathbf{x}^{(t)}$$

- Must keep track of counts of each feature (weighed by $y^{(t)}$):
 - E.g., with sparse data, for each non-zero dimension i in $\mathbf{x}^{(t)}$:

- Can we generalize the Bloom Filter?

©Carlos Guestrin 2013

17

Count-Min Sketch: single vector

- Simpler problem: Count how many times you see each string
- Single hash function:
 - Keep Count vector of length m
 - every time see string i :

$$\text{Count}[h(i)] \leftarrow \text{Count}[h(i)] + 1$$

- Again, collisions could be a problem:
 - a_i is the count of element i :

©Carlos Guestrin 2013

18

Count-Min Sketch: general case

- Keep d by m Count matrix
- d hash functions:
 - Just like in Bloom Filter, decrease errors with multiple hashes
 - Every time see string i :

$$\forall j \in \{1, \dots, d\} : \text{Count}[j, h(i)] \leftarrow \text{Count}[j, h(i)] + 1$$

Querying the Count-Min Sketch

$$\forall j \in \{1, \dots, d\} : \text{Count}[j, h(i)] \leftarrow \text{Count}[j, h(i)] + 1$$

- Query $Q(i)$?
 - What is in $\text{Count}[j, k]$?
 - Thus:
 - Return:

Analysis of Count-Min Sketch

$$\hat{a}_i = \min_j \text{Count}[j, h(i)] \geq a_i$$

- Set:
$$m = \left\lceil \frac{e}{\epsilon} \right\rceil \quad d = \left\lceil \ln \frac{1}{\delta} \right\rceil$$

- Then, after seeing n elements:

$$\hat{a}_i \leq a_i + \epsilon n$$

- With probability at least $1-\delta$

©Carlos Guestrin 2013

21

Proof of Count-Min for Point Query with Positive Counts: Part 1 – Expected Bound

- $I_{i,j,k}$ = indicator that i & k collide on hash j :
- Bounding expected value:
- $X_{i,j}$ = total colliding mass on estimate of count of i in hash j :
- Bounding colliding mass:
- Thus, estimate from each hash function is close in expectation

©Carlos Guestrin 2013

22

Proof of Count-Min for Point Query with Positive Counts: Part 2 – High Probability Bounds

- What we know: $Count[j, h(i)] = a_i + X_{i,j}$ $E[X_{i,j}] \leq \frac{\epsilon}{e} n$

- Markov inequality: For z_1, \dots, z_k positive iid random variables

$$P(\forall z_i : z_i > \alpha E[z_i]) < \alpha^{-k}$$

- Applying to the Count-Min sketch:

©Carlos Guestrin 2013

23

But Our Updates may be positive or Negative

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right] y^{(t)} \mathbf{x}^{(t)}$$

- Count-Min sketch for positive & negative case

- a_i no longer necessarily positive

- Update the same: Observe change Δ_i to element i :

$$\forall j \in \{1, \dots, d\} : Count[j, h(i)] \leftarrow Count[j, h(i)] + \Delta_i$$

- Each $Count[j, h(i)]$ no longer an upper bound on a_i

- How do we make a prediction?

- Bound: $|\hat{a}_i - a_i| \leq 3\epsilon \|\mathbf{a}\|_1$

- With probability at least $1 - \delta^{1/4}$, where $\|\mathbf{a}\| = \sum_i |a_i|$

©Carlos Guestrin 2013

24

Finally, Sketching for Perceptron

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \mathbb{1} \left[y^{(t)} (\mathbf{w}^{(t)} \cdot \mathbf{x}^{(t)}) \leq 0 \right] y^{(t)} \mathbf{x}^{(t)}$$

- Never need to know size of vocabulary!
- Make a mistake, update Count-Min matrix:

- Making a prediction:

- Scales to huge problems, great practical implications... More next time

©Carlos Guestrin 2013

25

What you need to know

- Update rules for perceptron and SGD problematic in high dimensions:
 - Complexity linear in number of features
 - What if we have a new feature?
- Hash functions
- Bloom filter
 - Test membership with some false positives, but very small number of bits per element
- Count-Min sketch
 - Positive counts: upper bound with nice rates of convergence
 - General case
- Application to Perceptron Learning and Prediction

©Carlos Guestrin 2013

26