

SCALABLE L1-REGULARIZED LOGISTIC REGRESSION:

IRLS-LARS

ANDREA KLEIN & DIPAN KUMAR PAL

THE PROBLEM

- ★ WE WISH TO DO BINARY CLASSIFICATION ON LARGE DATASETS VIA **LOGISTIC REGRESSION** WITH AN L1 NORM.
- ★ RECALL THAT LOGISTIC REGRESSION MODELS THE PROBABILITY THAT AN ELEMENT x BELONGS TO A CLASS y AS A SIGMOID FUNCTION:

$$p(y = 1|x, \theta) = \sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}$$

THE PROBLEM

- ★ UNDER A LAPLACIAN PRIOR, THE MAP ESTIMATE FOR THE MODEL PARAMETERS (THETA) IS GIVEN BY:

$$\min_{\theta} \sum_{i=1}^M -\log p(y^{(i)} | \mathbf{x}^{(i)}; \theta) + \beta \|\theta\|_1$$

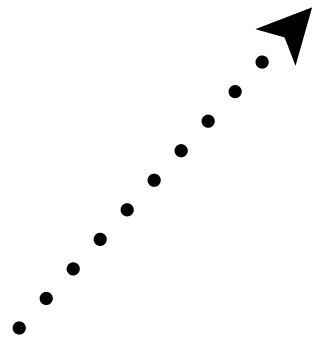
- ★ TECHNICALLY WE SOLVE AN EQUIVALENT ALTERNATIVE PARAMETRIZATION, SUBJECT TO:

$$\|\theta\|_1 \leq C$$

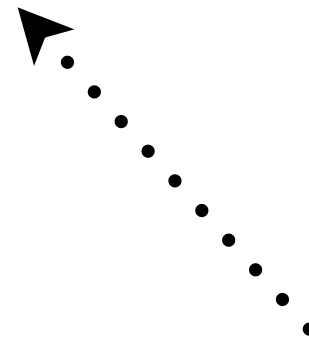
- ★ REGARDLESS OF HOW YOU FORMULATE THE PROBLEM, IT'S ORDINARILY INEFFICIENT TO FIND THETA FOR LARGE DATASETS.

THE ALGORITHM

IRLS-LARS



**ITERATIVELY RE-WEIGHTED
LEAST SQUARES**



**LEAST ANGLE
REGRESSION**

THE ALGORITHM

IRLS-LARS



ITERATIVELY RE-WEIGHTED
LEAST SQUARES

LEAST ANGLE
REGRESSION

- ★ ON EACH STEP OF THE ALGORITHM, WE CHARACTERIZE THE UPDATE DIRECTION (GAMMA) OF THETA AS THE SOLUTION TO A **LEAST-SQUARES PROBLEM**:

$$\gamma^{(k)} = \arg \min_{\gamma} \|(\Lambda^{\frac{1}{2}} \mathbf{X}^{\top})\gamma - \Lambda^{\frac{1}{2}} \mathbf{z}\|_2^2$$

THE ALGORITHM

IRLS-LARS

$$\mathbf{X} = [\mathbf{x}^{(1)} \mathbf{x}^{(2)} \dots \mathbf{x}^{(M)}]$$

UPDATED
("RE-WEIGHTED")
ON EACH ITERATION:

$$\begin{aligned}\Lambda_{ii} &= \sigma(\theta^{(k)\top} \mathbf{x}^{(i)}) [1 - \sigma(\theta^{(k)\top} \mathbf{x}^{(i)})], \\ z_i &= \mathbf{x}^{(i)\top} \theta^{(k)} + \frac{[1 - \sigma(y^{(i)} \theta^{(k)\top} \mathbf{x}^{(i)})] y^{(i)}}{\Lambda_{ii}}.\end{aligned}$$

ITERATIVELY

LEAST SQUARES

ANGLE

REGRESSION

- ★ ON EACH STEP OF THE ALGORITHM, WE CHARACTERIZE THE UPDATE DIRECTION (GAMMA) OF THETA AS THE SOLUTION TO A **LEAST-SQUARES PROBLEM**:

$$\gamma^{(k)} = \arg \min_{\gamma} \|(\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{X}^{\top}) \gamma - \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{z}\|_2^2$$

THE ALGORITHM

IRLS-LARS



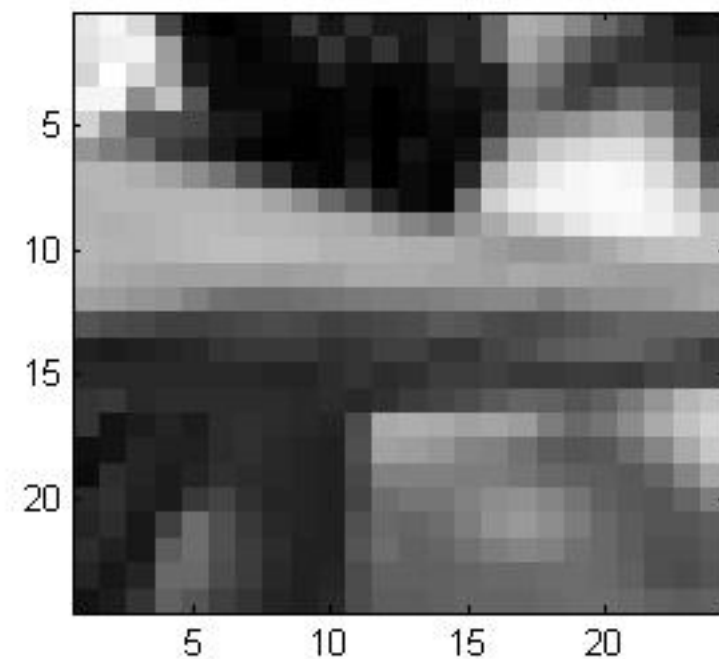
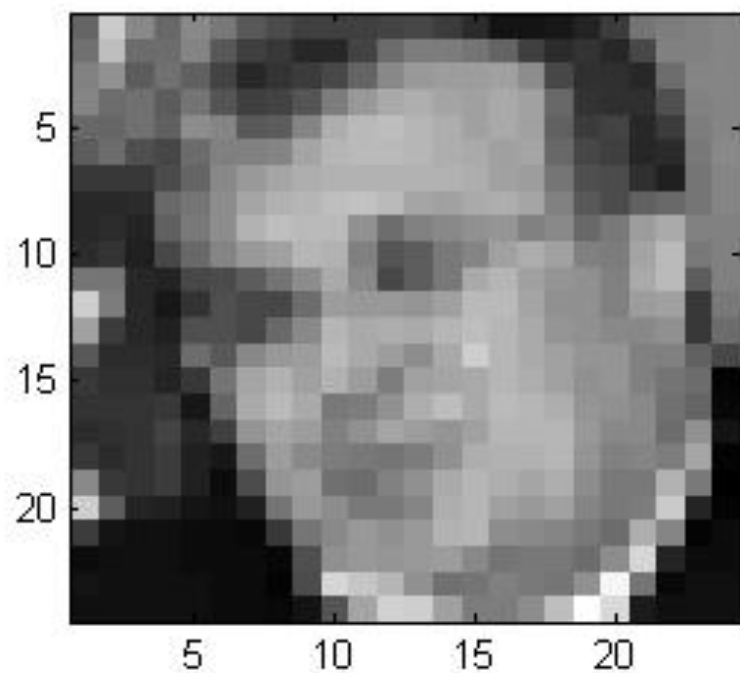
ITERATIVELY RE-WEIGHTED
LEAST SQUARES

LEAST ANGLE
REGRESSION

- ★ IN PRACTICE, WE USE AN ALGORITHM CALLED **LARS** TO EFFICIENTLY SOLVE FOR THE STEP DIRECTION WITHIN THE IRLS FORMULATION. FINALLY, ONCE WE KNOW THE DIRECTION OF THE UPDATE, WE USE A BACKTRACKING LINE SEARCH TO DETERMINE HOW BIG A STEP TO TAKE.

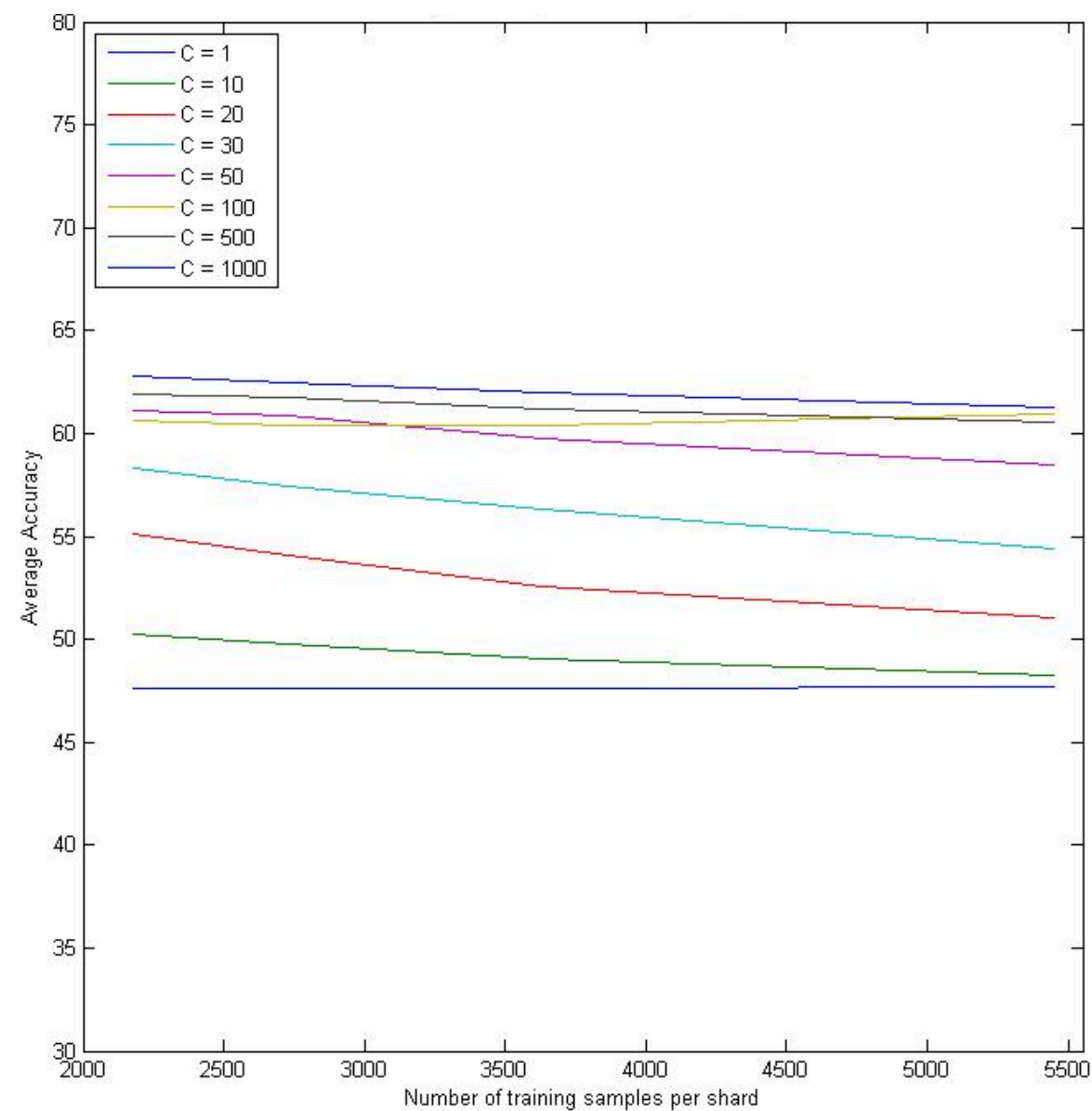
NOW FOR AN EXPERIMENT

WE'LL DEMONSTRATE IRLS-LARS ON A BINARY IMAGE CLASSIFICATION TASK: IDENTIFYING WHICH IMAGES HAVE FACES (AND WHICH DON'T).



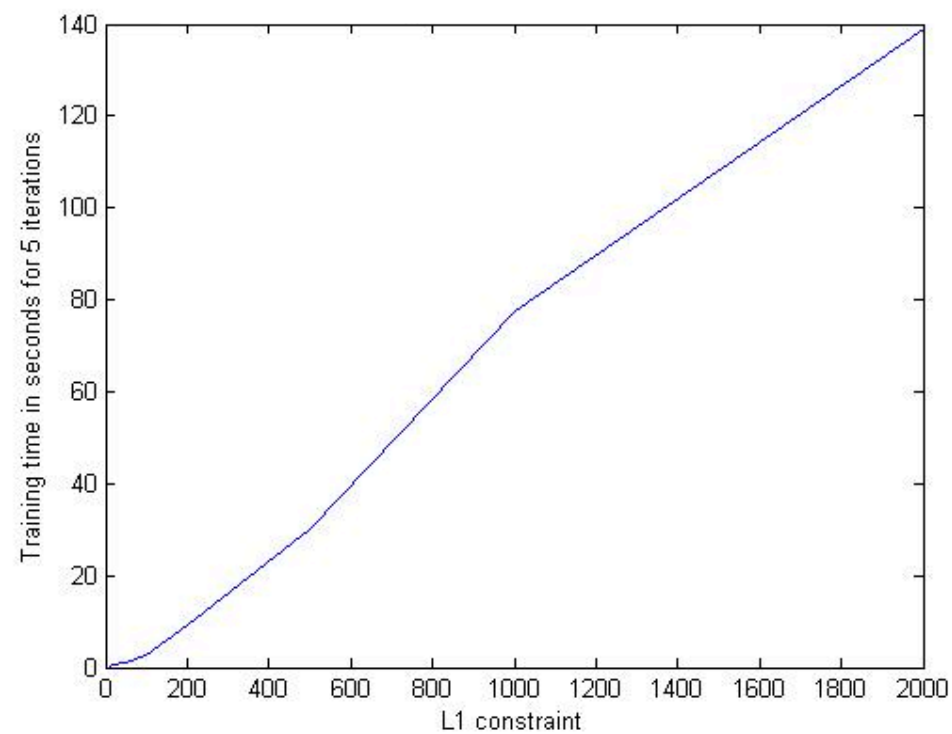
PERFORMANCE

AVERAGE ACCURACY VS. # OF TRAINING SAMPLES

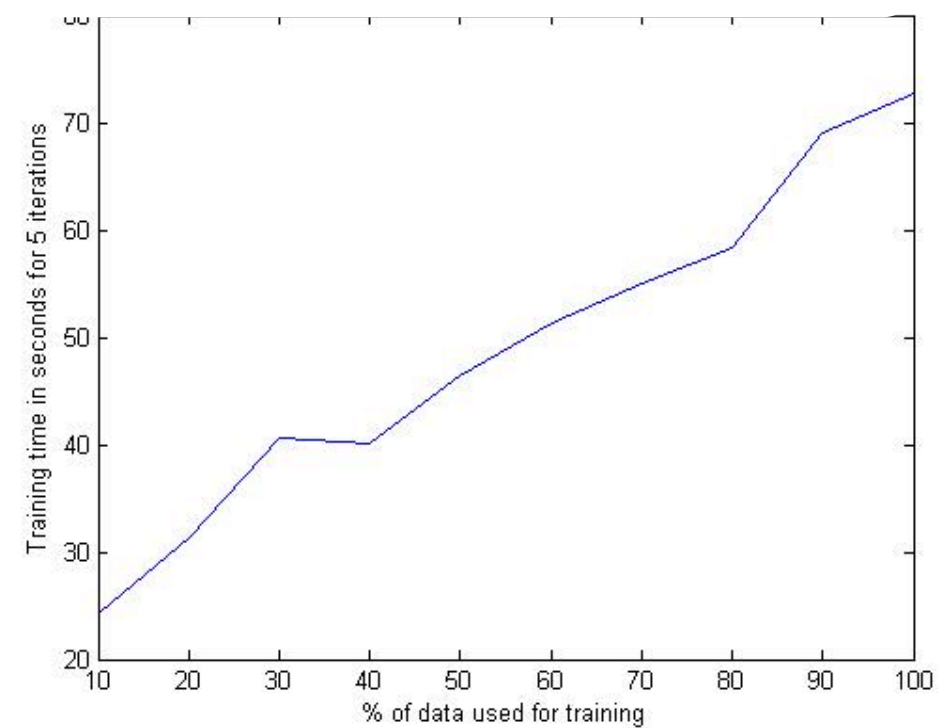


PERFORMANCE

TRAINING TIME VS. CONSTRAINT SIZE



TRAINING TIME VS. TRAINING SIZE



SCALABILITY

Algorithm 2 ParallelSGD($\{c^1, \dots, c^m\}, T, \eta, w_0, k$)

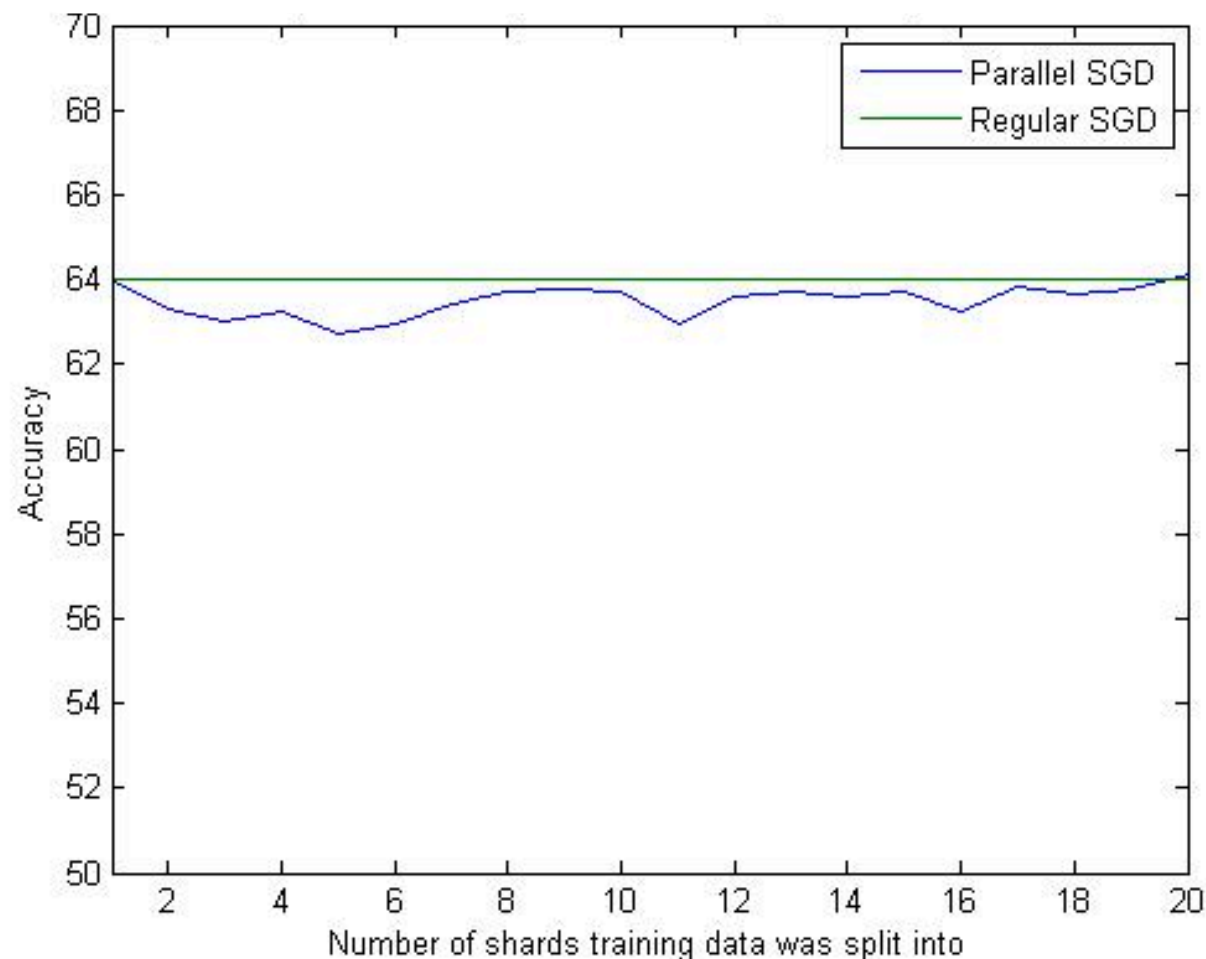
for all $i \in \{1, \dots, k\}$ **parallel do**

$v_i = \text{SGD}(\{c^1, \dots, c^m\}, T, \eta, w_0)$ on client

end for

Aggregate from all computers $v = \frac{1}{k} \sum_{i=1}^k v_i$ and **return** v

RELATIVE PERFORMANCE OF PARALLEL VS. REGULAR SGD



STILL TO COME

- ★ MORE STRENUOUS (PARALLEL) TESTS ON LARGER DATASETS
- ★ IF TIME, WE'D LIKE TO EXPLORE THE ALGORITHM'S PERFORMANCE ON MORE DIVERSE DATASETS
- ★ WE MAY ALSO INCORPORATE **PARALLEL SCD** (SEQUENTIAL COORDINATE DESCENT), PARALLELIZING OVER FEATURES RATHER THAN SAMPLES:

Algorithm 2 Shotgun: Parallel SCD

Choose number of parallel updates $P \geq 1$.

Set $\mathbf{x} = \mathbf{0} \in \mathbb{R}_+^{2d}$

while not converged **do**

 Choose random subset of P weights in $\{1, \dots, 2d\}$.

In parallel on P processors

 Get assigned weight j .

 Set $\delta x_j \leftarrow \max\{-x_j, -(\nabla F(\mathbf{x}))_j/\beta\}$.

 Update $x_j \leftarrow x_j + \delta x_j$.

end while

THE END!

REFERENCES:

- ★ EFFICIENT L1 REGULARIZED LOGISTIC REGRESSION (LEE ET AL)
- ★ PARALLELIZED STOCHASTIC GRADIENT DESCENT (ZINKEVICH ET AL)
- ★ PARALLEL COORDINATE DESCENT FOR L1 REGULARIZED LOSS MINIMIZATION (BRADLEY ET AL)
- ★ LEAST ANGLE REGRESSION (EFRON ET AL)

MATH APPENDIX

$$\mathbf{X} = [\mathbf{x}^{(1)} \mathbf{x}^{(2)} \dots \mathbf{x}^{(M)}].$$

Let the diagonal matrix $\mathbf{\Lambda}$ and the vector \mathbf{z} be defined as follows: for all $i = 1, 2, \dots, M$:

$$\Lambda_{ii} = \sigma(\theta^{(k)\top} \mathbf{x}^{(i)}) [1 - \sigma(\theta^{(k)\top} \mathbf{x}^{(i)})], \quad (7)$$

$$z_i = \mathbf{x}^{(i)\top} \theta^{(k)} + \frac{[1 - \sigma(y^{(i)} \theta^{(k)\top} \mathbf{x}^{(i)})] y^{(i)}}{\Lambda_{ii}}. \quad (8)$$

Then, we have that $\mathbf{H}(\theta^{(k)}) = -\mathbf{X}\mathbf{\Lambda}\mathbf{X}^\top$ and $\mathbf{g}(\theta^{(k)}) = \mathbf{X}\mathbf{\Lambda}(\mathbf{z} - \mathbf{X}^\top \theta^{(k)})$, and thus Equation (5) can be rewritten as:

$$\gamma^{(k)} = (\mathbf{X}\mathbf{\Lambda}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{\Lambda}\mathbf{z}. \quad (9)$$

Thus $\gamma^{(k)}$ is the solution to the following weighted least squares problem:

$$\gamma^{(k)} = \arg \min_{\gamma} \|(\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{X}^\top) \gamma - \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{z}\|_2^2. \quad (10)$$