

Problem

Efficiently solve nonnegative Lasso:

$$\begin{aligned} & \text{minimize}_{\mathbf{x} \in \mathbb{R}^m} && \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 + \lambda \|\mathbf{x}\|_1 \\ & \text{s.t.} && \mathbf{x} \geq \mathbf{0} \end{aligned}$$

Solution is *sparse*—most weights have value 0

StingyCD can also solve:

- Unconstrained Lasso
- Linear SVM

Wasteful Updates in Coordinate Descent

Coordinate descent is a popular and efficient algorithm for solving this problem

Coordinate descent

```
for t = 1, 2, ... do
  i ← get_next_coordinate()
  δ ← compute_update(i, x^(t-1))
  x^(t) ← x^(t-1) + δe_i
```

Define residuals vector

$$\mathbf{r}^{(t)} = \mathbf{b} - \mathbf{Ax}^{(t)}$$

During iteration t , choose i and compute

$$g_i \leftarrow -\langle \mathbf{A}_i, \mathbf{r}^{(t-1)} \rangle + \lambda \quad \delta \leftarrow \max \left\{ -x_i^{(t-1)}, \frac{-g_i}{\|\mathbf{A}_i\|^2} \right\}$$

“Zero updates” During iteration t , $\delta = 0$ if

1. $g_i \geq 0$ and
2. $x_i^{(t-1)} = 0$

Due to sparsity, often *majority* of updates are zero

Zero updates are *wasteful*—computing g_i requires time

Geometry of zero update

$$\bullet \mathbf{r}^{(t-1)} = \mathbf{b} - \mathbf{Ax}^{(t-1)}$$

$$\mathbf{r}^{(t-1)} \notin \mathcal{A}_i \Leftrightarrow g_i \geq 0$$

$$\mathcal{A}_i = \{\mathbf{r} : -\langle \mathbf{A}_i, \mathbf{r} \rangle + \lambda < 0\}$$

StingyCD

StingyCD skips over many updates guaranteed to be zero

Testing skip condition requires negligible time

Reference residuals (updated infrequently)

$$\mathbf{rr} = \mathbf{r}^{(t-k)}$$

Reference distance (updated after each nonzero update)

$$d^{(t)} = \|\mathbf{r}^{(t)} - \mathbf{rr}\|_2$$

Update to reference distance requires negligible time

$$(d^{(t)})^2 = \underbrace{(d^{(t-1)})^2}_{\text{Cached quantities}} - 2\delta \underbrace{\langle \mathbf{A}_i, \mathbf{r}^{(t-1)} \rangle}_{\text{Cached quantities}} + 2\delta \underbrace{\langle \mathbf{A}_i, \mathbf{rr} \rangle}_{\text{Cached quantities}} + \delta^2 \underbrace{\|\mathbf{A}_i\|^2}_{\text{Cached quantities}}$$

Active region distances (computed after each \mathbf{rr} update)

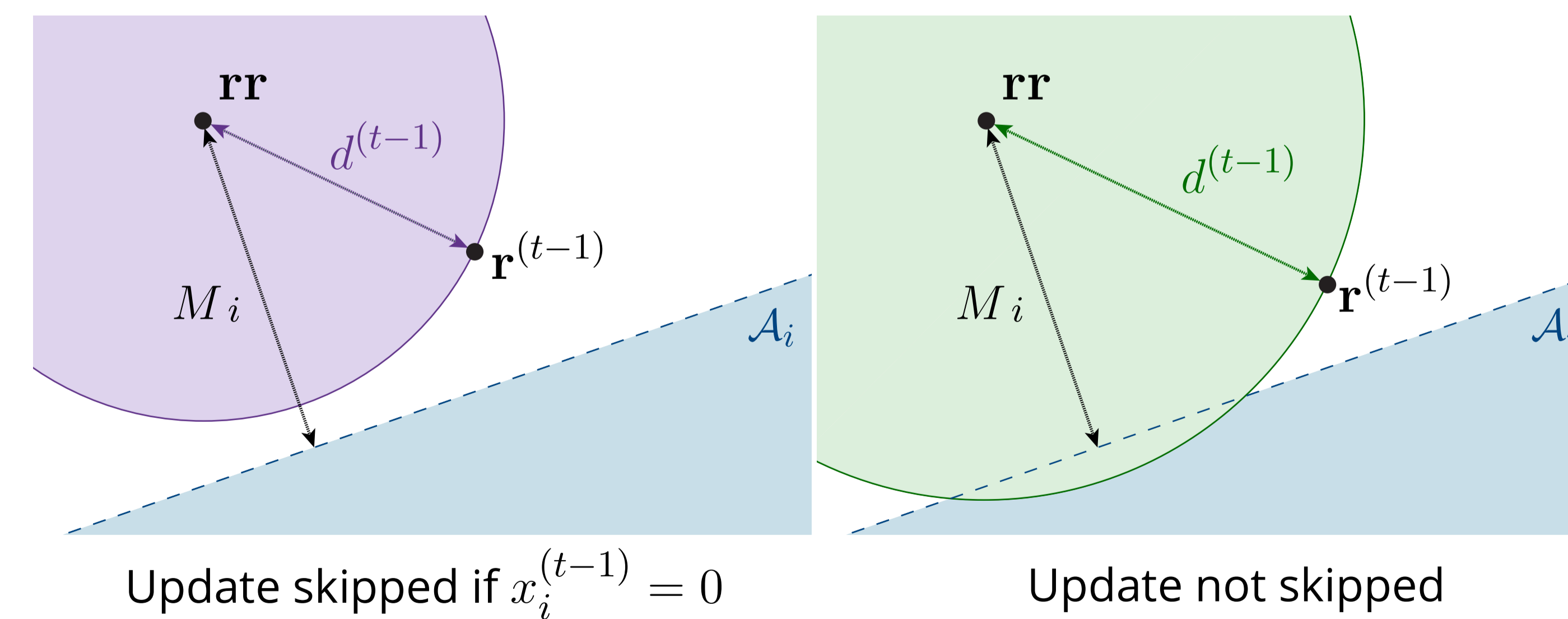
$$\text{For all } i, \quad M_i = \frac{1}{\|\mathbf{A}_i\|} (\lambda - \langle \mathbf{A}_i, \mathbf{rr} \rangle)$$

Distance between \mathbf{rr} and \mathcal{A}_i (but negative if $\mathbf{rr} \in \mathcal{A}_i$)

Full pass over data to compute M_i for all i

Skip condition StingyCD skips update t if

1. $d^{(t-1)} \leq M_i$ and
2. $x_i^{(t-1)} = 0$



Safeness guarantee $d^{(t-1)} \leq M_i \Rightarrow \mathbf{r}^{(t-1)} \notin \mathcal{A}_i \Rightarrow g_i \geq 0$

In StingyCD, every skipped update would, if computed, result in $\delta = 0$

StingyCD

```
for t = 1, 2, ... do
  if update_reference?() then
    rr ← r^(t-1)
    d^(t-1) ← 0
    Mi ← compute_M(rr) for all i

  i ← get_next_coordinate()

  if d^(t-1) ≤ Mi and xi^(t-1) = 0 then
    x^(t) ← x^(t-1); r^(t) ← r^(t-1); ...
    continue

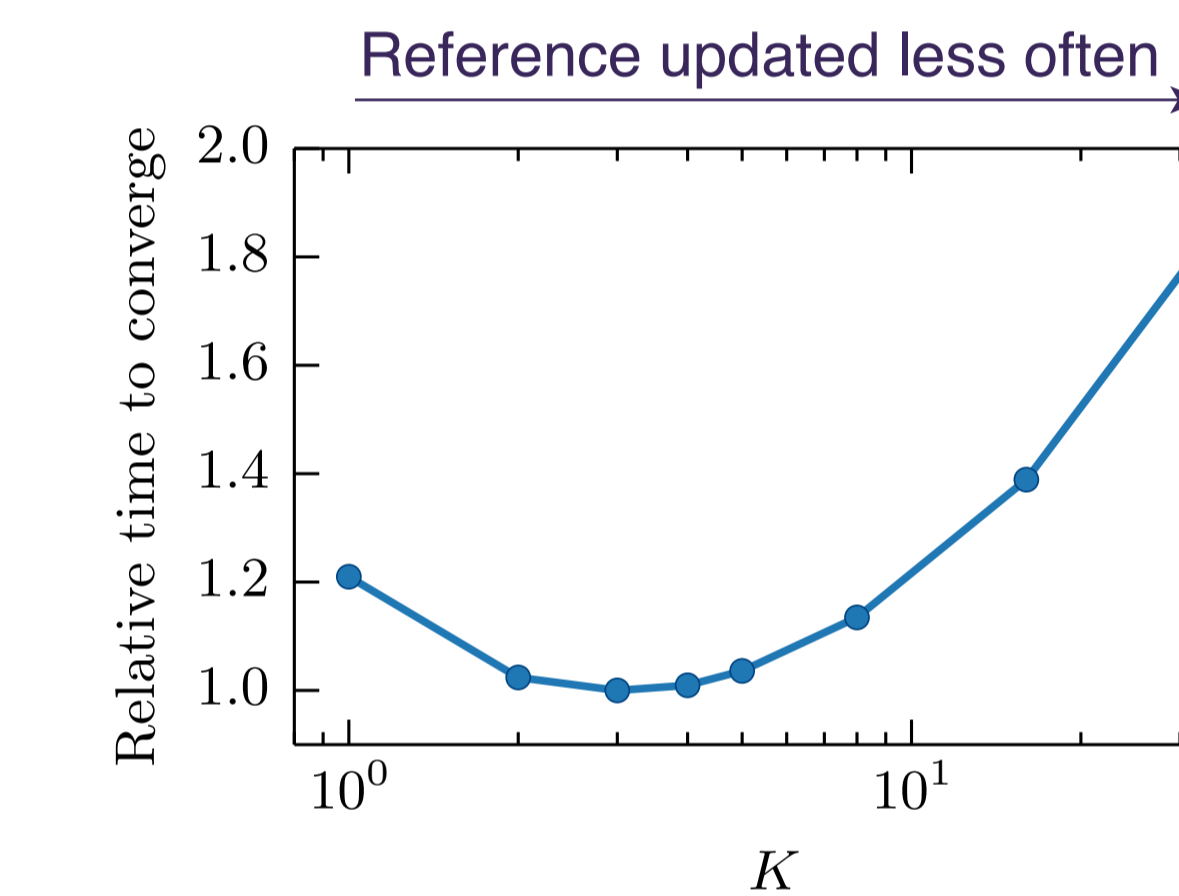
  δ ← compute_update(i, x^(t-1))
  x^(t) ← x^(t-1) + δe_i
  r^(t) ← r^(t-1) - δAi
  d^(t) ← update_d(d^(t-1), δ, i)
```

Compared to CD, a StingyCD iteration uses...

- **Less time** if update is skipped
- **More time** if reference is updated
- **The same amount of time** otherwise

Scheduling reference updates

Measure time T to first update \mathbf{rr}
Update \mathbf{rr} after blocks of KT time
We set $K = 5$



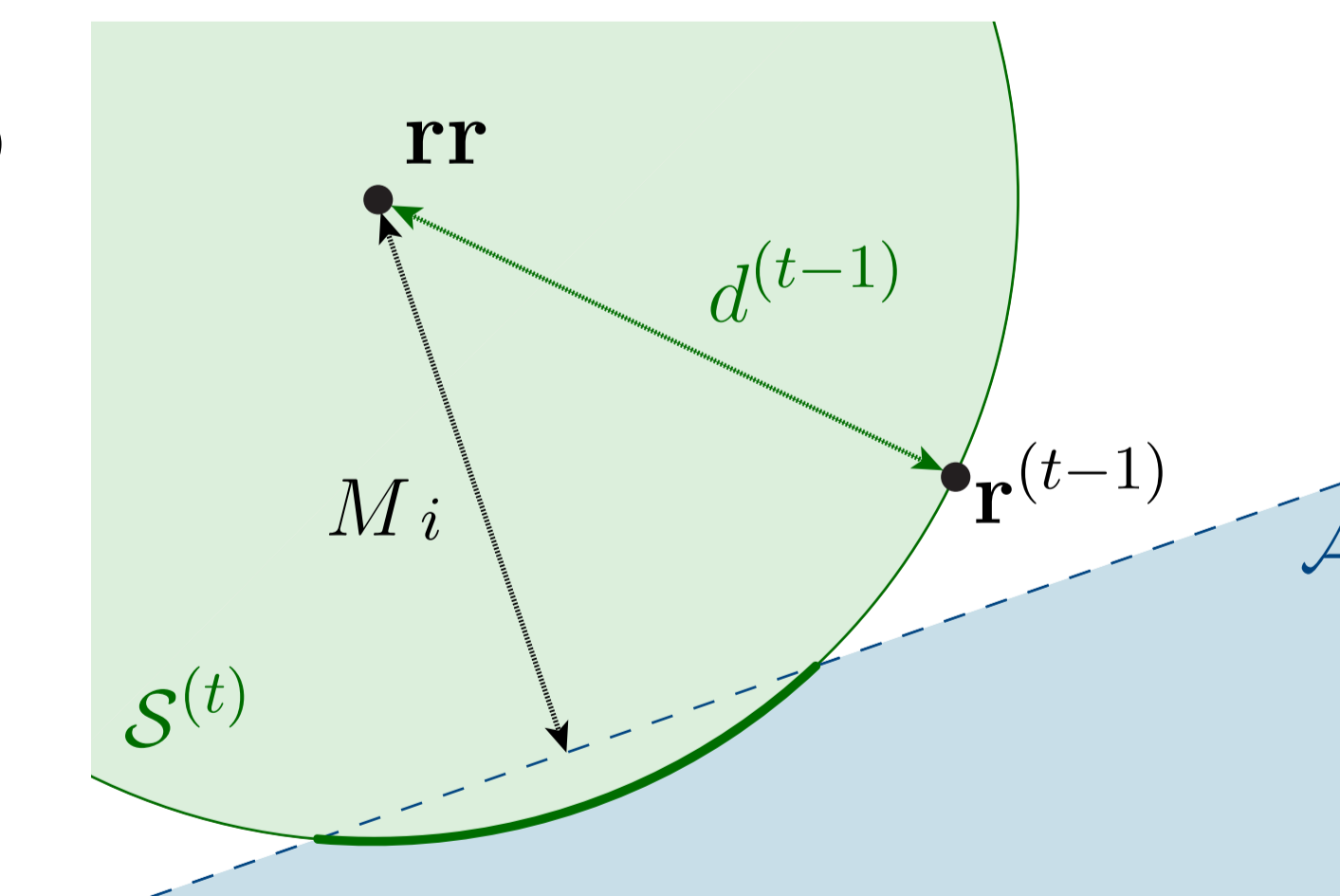
StingyCD+

StingyCD-inspired heuristic—skips updates more aggressively

Models probability of nonzero update, denoted $P(U^{(t)})$

If $x_i^{(t-1)} = 0$, then

$$P(U^{(t)}) = \frac{\text{Area}(\text{bd}(\mathcal{S}^{(t)}) \cap \mathcal{A}_i)}{\text{Area}(\text{bd}(\mathcal{S}^{(t)}))}$$



Skip update if “expected relevant delay” is small:

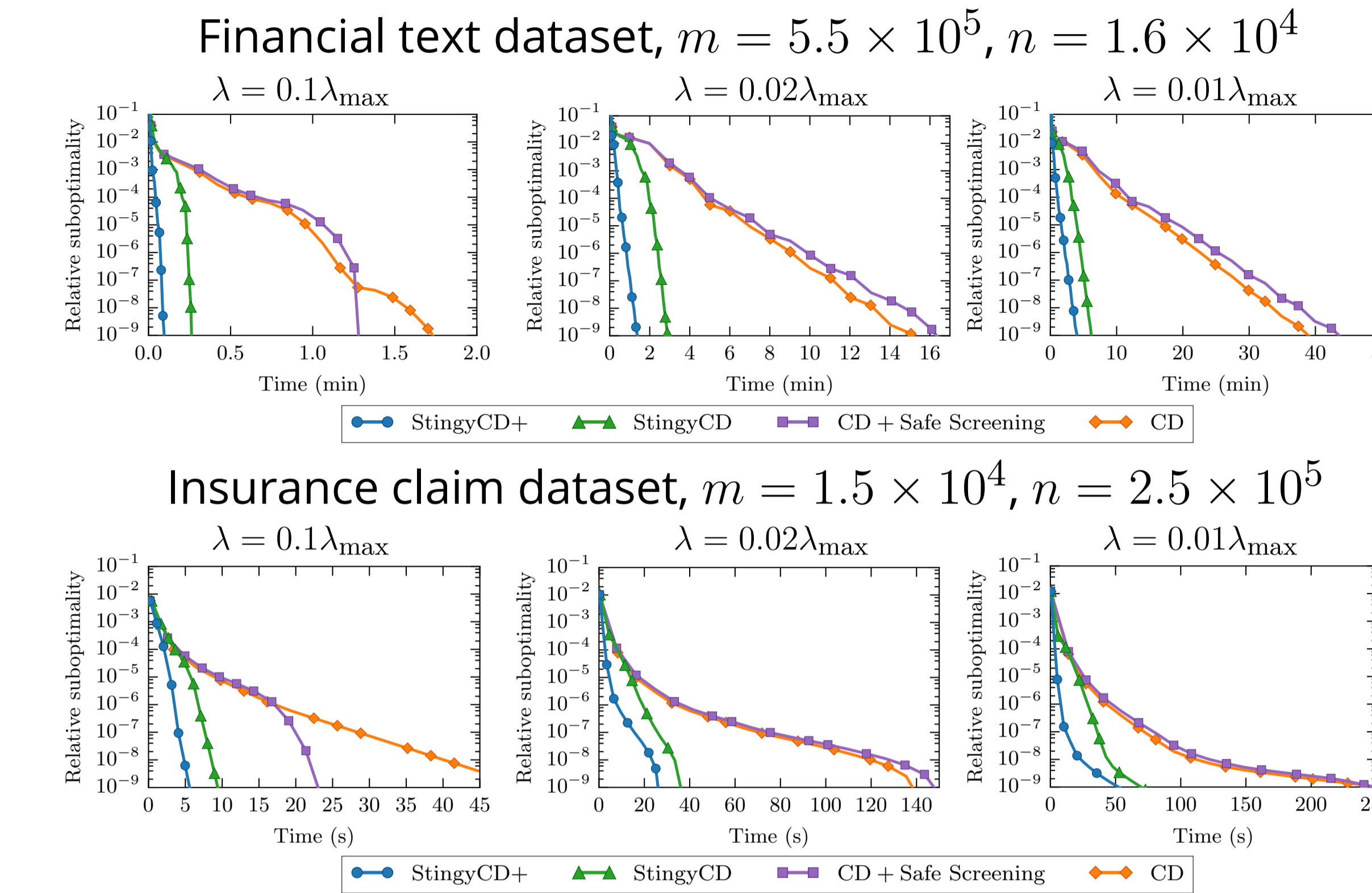
$$P(U^{(t)}) D_i^{(t)} < \xi^{(t)}$$

Number of updates computed since coordinate i last updated

We set $\xi^{(t)} = \text{NNZ}(\mathbf{x}^{(t-1)})$ (still converges to solution)

Empirical Results

Lasso problems

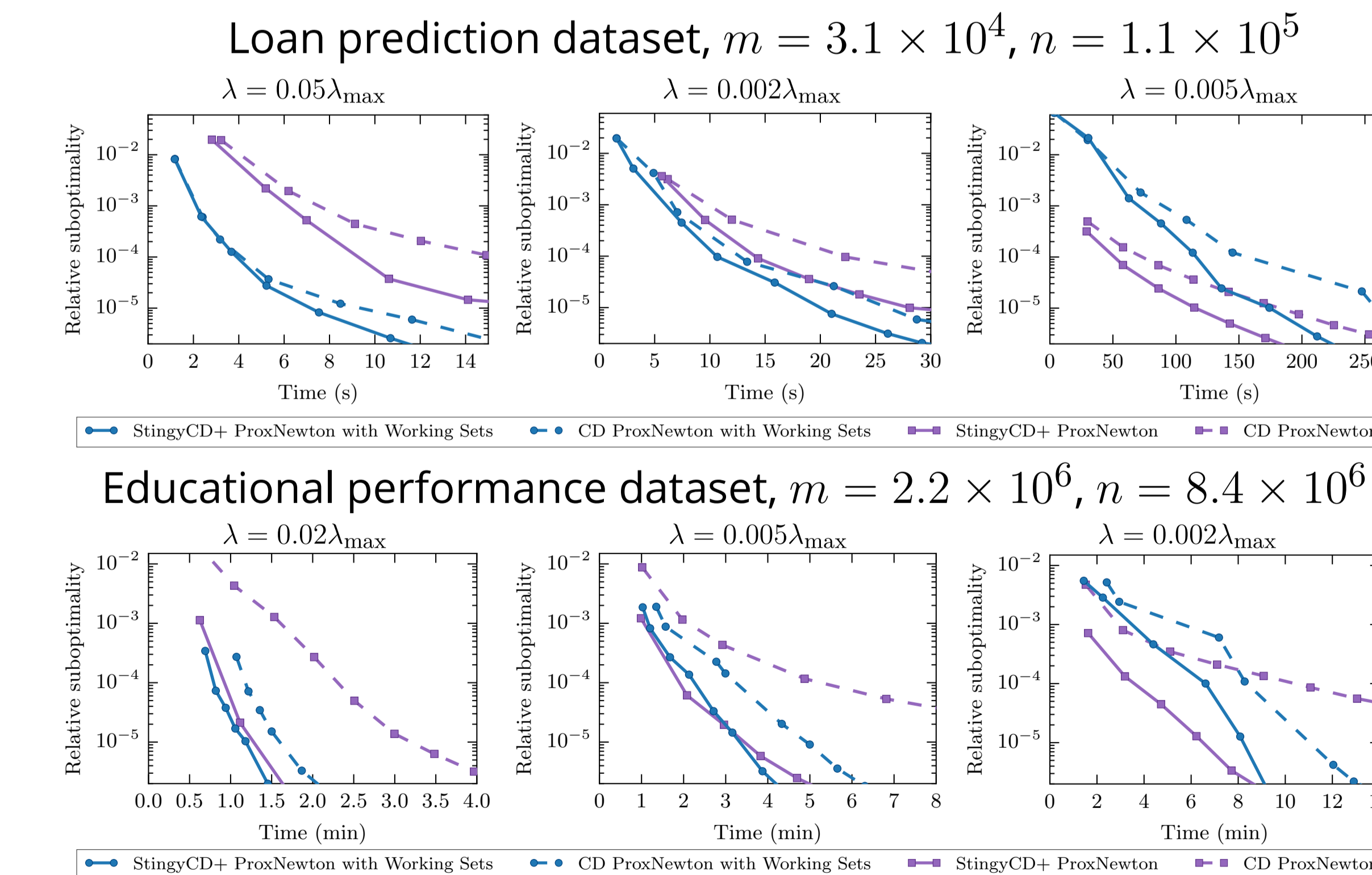


Sparse logistic regression problems

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^m} \sum_{i=1}^n \log(1 + \exp(-b_i \mathbf{a}_i^T \mathbf{x})) + \lambda \|\mathbf{x}\|_1$$

Proximal Newton method is efficient—solves sequence of Lasso problems

Compare CD, StingyCD+ as subproblem solvers



Conclusions

StingyCD can avoid many wasteful CD updates
Relaxations (StingyCD+) can lead to further gains
Can also combine StingyCD with other methods