
Parallel Coordinate Descent for L_1 -Regularized Loss Minimization: Theory Supplement

Abstract

In this supplementary document, we give detailed proofs of all theoretical results of the main paper.

1. Preliminaries

General form for our optimization problems, modified to use duplicate features and have a twice-differentiable regularization term:

$$\min_{\mathbf{x} \in \mathbb{R}_+^{2d}} \sum_{i=1}^n L(\hat{\mathbf{a}}_i^T \mathbf{x}, y_i) + \lambda \sum_{j=1}^{2d} x_j \quad (1)$$

Instantiation of Eq. (1) for Lasso (?):

$$F(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^{2d} x_j \quad (2)$$

Instantiation of Eq. (1) for sparse logistic regression:

$$F(\mathbf{x}) = \sum_{i=1}^n \log(1 + \exp(-y_i \hat{\mathbf{a}}_i^T \mathbf{x})) + \lambda \sum_{j=1}^{2d} x_j \quad (3)$$

Update rule for $x_j \leftarrow x_j + \delta x_j$:

$$\delta x_j = \max\{-x_j, -(\nabla F(\mathbf{x}))_j / \beta\} \quad (4)$$

2. Detailed Proofs: β for Squared Loss and Logistic Loss

Assumptions 2.1 and 3.1 both upper bound the change in objective from updating \mathbf{x} with $\Delta \mathbf{x}$. We show how to do so for Assumption 3.1, which generalizes Assumption 2.1. For both losses, we upper-bound the objective using a second-order Taylor expansion of F around \mathbf{x} .

Appearing in *Proceedings of the 28th International Conference on Machine Learning*, Bellevue, WA, USA, 2011. Copyright 2011 by the author(s)/owner(s).

2.1. Proof: β for Squared Loss

$$\begin{aligned} \nabla F(\mathbf{x}) &= \mathbf{A}^T \mathbf{A} \mathbf{x} - \mathbf{A}^T \mathbf{y} + \lambda \mathbf{1} & (5) \\ \nabla^2 F(\mathbf{x}) &= \mathbf{A}^T \mathbf{A} & (6) \end{aligned}$$

where $\mathbf{1}$ is an all-ones vector of the appropriate size. Note that, since derivatives of (2) of order higher than two are zero, the second order Taylor expansion is exact:

$$\begin{aligned} F(\mathbf{x} + \Delta \mathbf{x}) &= F(\mathbf{x}) + (\Delta \mathbf{x})^T \nabla F(\mathbf{x}) + \frac{1}{2} (\Delta \mathbf{x})^T \nabla^2 F(\mathbf{x}) \Delta \mathbf{x} & (7) \end{aligned}$$

Plugging in the second order derivative gives $\beta = 1$:

$$\begin{aligned} F(\mathbf{x} + \Delta \mathbf{x}) &= F(\mathbf{x}) + (\Delta \mathbf{x})^T \nabla F(\mathbf{x}) + \frac{1}{2} (\Delta \mathbf{x})^T \mathbf{A}^T \mathbf{A} \Delta \mathbf{x}. & (8) \end{aligned}$$

This bound is exact for squared loss but not for all losses.

2.2. Proof: β for Logistic Loss

Define $p_i = \frac{1}{1 + \exp(-\mathbf{a}_i^T \mathbf{x})}$, the class conditional probability of y_i given \mathbf{a}_i .

$$\frac{\partial}{\partial x_j} F(\mathbf{x}) = \lambda + \sum_{i=1}^n y_i \mathbf{A}_{ij} (p_i - 1) \quad (9)$$

$$\frac{\partial^2}{\partial x_j \partial x_k} F(\mathbf{x}) = \sum_{i=1}^n \mathbf{A}_{ij} \mathbf{A}_{ik} (1 - p_i) p_i \quad (10)$$

Taylor's theorem tells us that there exists an $\hat{\mathbf{x}}$ s.t.

$$\begin{aligned} F(\mathbf{x} + \Delta \mathbf{x}) &\leq F(\mathbf{x}) + (\Delta \mathbf{x})^T \nabla F(\mathbf{x}) + \frac{1}{2} (\Delta \mathbf{x})^T (\nabla^2 F(\hat{\mathbf{x}})) \Delta \mathbf{x} & (11) \end{aligned}$$

The second-order term is maximized by setting $p_i = \frac{1}{2}$ in $\frac{\partial^2}{\partial x_j \partial x_k} F(\mathbf{x})$ for each j, k . Plugging this in gives our bound with $\beta = \frac{1}{4}$:

$$\begin{aligned} F(\mathbf{x} + \Delta \mathbf{x}) &\leq F(\mathbf{x}) + (\Delta \mathbf{x})^T \nabla F(\mathbf{x}) + \frac{1}{2} \frac{(\Delta \mathbf{x})^T \mathbf{A}^T \mathbf{A} \Delta \mathbf{x}}{4} & (12) \end{aligned}$$

3. Duplicated Features

Our work, like Shalev-Shwartz and Tewari (?)’s work, uses duplicated features (with $\mathbf{x} \in \mathbb{R}^{2d}$ and $\mathbf{A} \in \mathbb{R}^{n \times 2d}$), but our actual algorithm does not (so $\mathbf{x} \in \mathbb{R}_+^d$ and $\mathbf{A} \in \mathbb{R}^{n \times d}$). They point out that the optimization problems with and without duplicated features are equivalent.

To see this, consider the form of $F(\mathbf{x})$ in Eq. (1). x_j only appears in the dot product $\hat{\mathbf{a}}_i^T \mathbf{x}$ via $\mathbf{A}_{i,j} x_j$, and x_{d+j} only appears via $-\mathbf{A}_{i,j} x_{d+j}$, where \mathbf{A} is the original design matrix without duplicated features. Suppose $x_j > 0$ and $x_{d+j} > 0$, and assume w.l.o.g. that $x_j > x_{d+j}$. Then setting $x_j \leftarrow x_j - x_{d+j}$ and $x_{d+j} \leftarrow 0$ would give the same value for the loss term $L(\hat{\mathbf{a}}_i^T \mathbf{x}, y_i)$, and it would decrease the regularization penalty by $2\lambda x_{d+j}$. Therefore, at the optimum, at most one of x_j, x_{d+j} will be non-zero, and the objectives with and without duplicated features will be equal.

4. Detailed Proof: Theorem 2.1

Define a potential function, where \mathbf{x}^* is the optimal weight vector:

$$\Psi(\mathbf{x}) = \frac{\beta}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 + F(\mathbf{x}) \quad (13)$$

Claim: After updating weight x_j with δx_j ,

$$\Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x}) \geq (x_j - x_j^*)(\nabla F(\mathbf{x}))_j \quad (14)$$

To see this:

$$\begin{aligned} \Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x}) &= \frac{\beta}{2} [\|\mathbf{x} - \mathbf{x}^*\|_2^2 - \|\mathbf{x} + \Delta \mathbf{x} - \mathbf{x}^*\|_2^2] \\ &\quad + F(\mathbf{x}) - F(\mathbf{x} + \Delta \mathbf{x}) \end{aligned} \quad (15)$$

$$\begin{aligned} &= -\frac{\beta}{2} [2(\mathbf{x} - \mathbf{x}^*)^T \Delta \mathbf{x} + (\delta x_j)^2] \\ &\quad + F(\mathbf{x}) - F(\mathbf{x} + \Delta \mathbf{x}) \end{aligned} \quad (16)$$

$$\geq \beta \left(-\mathbf{x}^T \Delta \mathbf{x} + \mathbf{x}^{*T} \Delta \mathbf{x} - \frac{(\delta x_j)^2}{2} \right) \quad (17)$$

$$\begin{aligned} &\quad - (\Delta \mathbf{x})^T \nabla F(\mathbf{x}) - \frac{\beta}{2} (\delta x_j)^2 \\ &= \beta (-x_j \delta x_j + x_j^* \delta x_j - (\delta x_j)^2) \end{aligned} \quad (18)$$

$$\begin{aligned} &\quad - \delta x_j (\nabla F(\mathbf{x}))_j \\ &\geq \beta (-x_j \delta x_j - (\delta x_j)^2) - x_j^* (\nabla F(\mathbf{x}))_j \\ &\quad - \delta x_j (\nabla F(\mathbf{x}))_j \end{aligned} \quad (19)$$

Above, Eq. (17) used Assumption 2.1. Eq. (19) used the update rule for choosing δx_j in Eq. (4). Now there are two possible cases stemming from the update rule. Case 1: If $\delta x_j = -x_j$, then Eq. (19) simplifies to

$$\Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x}) \geq (x_j - x_j^*)(\nabla F(\mathbf{x}))_j \quad (20)$$

Case 2: If $\delta x_j = -(\nabla F(\mathbf{x}))_j / \beta$, then Eq. (19) again simplifies to

$$\Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x}) \quad (21)$$

$$\begin{aligned} &\geq x_j (\nabla F(\mathbf{x}))_j - \beta (\delta x_j)^2 - x_j^* (\nabla F(\mathbf{x}))_j \\ &\quad + \beta (\delta x_j)^2 \end{aligned} \quad (22)$$

$$= (x_j - x_j^*) (\nabla F(\mathbf{x}))_j \quad (23)$$

Having proved our claim, we can now take the expectation of Eq. (14) w.r.t. j , the chosen weight:

$$\begin{aligned} \mathbf{E} [\Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x})] &\geq \mathbf{E} [(x_j - x_j^*) (\nabla F(\mathbf{x}))_j] \end{aligned} \quad (24)$$

$$= \frac{1}{2d} \mathbf{E} [(\mathbf{x} - \mathbf{x}^*)^T \nabla F(\mathbf{x})] \quad (25)$$

$$\geq \frac{1}{2d} \mathbf{E} [F(\mathbf{x}) - F(\mathbf{x}^*)] \quad (26)$$

In Eq. (25), we write $\frac{1}{2d}$ instead of $\frac{1}{d}$ (which Shalev-Shwartz and Tewari (?) write), for there are another d duplicates of each of the original d weights. Eq. (26) holds since $F(\mathbf{x})$ is convex.

Summing over $T + 1$ iterations gives:

$$\begin{aligned} \mathbf{E} \left[\sum_{t=0}^T \Psi(\mathbf{x}^{(t)}) - \Psi(\mathbf{x}^{(t+1)}) \right] &\geq \frac{1}{2d} \mathbf{E} \left[\sum_{t=0}^T F(\mathbf{x}^{(t)}) \right] - \frac{T+1}{2d} F(\mathbf{x}^*) \end{aligned} \quad (27)$$

$$\geq \frac{T+1}{2d} \left[\mathbf{E} [F(\mathbf{x}^{(T)})] - F(\mathbf{x}^*) \right] \quad (28)$$

where Eq. (28) used the fact that $F(\mathbf{x}_t)$ decreases monotonically with t . Since $\sum_{t=0}^T \Psi(\mathbf{x}^{(t)}) - \Psi(\mathbf{x}^{(t+1)}) = \Psi(\mathbf{x}^{(0)}) - \Psi(\mathbf{x}^{(T+1)})$, rearranging the above inequality gives

$$\mathbf{E} [F(\mathbf{x}_T)] - F(\mathbf{x}^*) \quad (29)$$

$$\leq \frac{2d}{T+1} \mathbf{E} [\Psi(\mathbf{x}^{(0)}) - \Psi(\mathbf{x}^{(T+1)})] \quad (30)$$

$$\leq \frac{2d}{T+1} \mathbf{E} [\Psi(\mathbf{x}^{(0)})] \quad (31)$$

$$= \frac{2d}{T+1} \left[\frac{\beta}{2} \|\mathbf{x}^*\|_2^2 + F(\mathbf{x}^{(0)}) \right] \quad (32)$$

where Eq. (31) used $\Psi(\mathbf{x}) \geq 0$ and Eq. (32) used $\mathbf{x}^{(0)} = 0$.

This bound divides by $T + 1$ instead of T (which Shalev-Shwartz and Tewari (?) do). Also, their theorem has an extra factor of $\frac{1}{2}$ on the right-hand side but should not due to the doubled length of \mathbf{x} (though careful analysis without duplicated features could likely re-introduce the $\frac{1}{2}$).

5. Detailed Proof: Theorem 3.1

Start with Eq. (8), and note that the update rule in Eq. (4) implies that $\delta x_j \leq -(\nabla F(\mathbf{x}))_j$ (with $\beta = 1$ for Lasso). This gives us:

$$\begin{aligned} F(\mathbf{x} + \Delta \mathbf{x}) - F(\mathbf{x}) \\ \leq -(\Delta \mathbf{x})^T (\Delta \mathbf{x}) + \frac{1}{2} (\Delta \mathbf{x})^T \mathbf{A}^T \mathbf{A} \Delta \mathbf{x} \end{aligned} \quad (33)$$

Noting that $\Delta \mathbf{x}$ can only have non-zeros in the indices in \mathcal{P}_t , we can rewrite this as

$$\begin{aligned} F(\mathbf{x} + \Delta \mathbf{x}) - F(\mathbf{x}) \\ \leq -\sum_{j \in \mathcal{P}_t} (\delta x_j)^2 + \frac{1}{2} \sum_{i,j \in \mathcal{P}_t} (\mathbf{A}^T \mathbf{A})_{i,j} \delta x_i \delta x_j \end{aligned} \quad (34)$$

Separating out the diagonal terms in the sum over i, j and using $\text{diag}(\mathbf{A}^T \mathbf{A}) = \mathbf{1}$ gives the desired result:

$$\begin{aligned} F(\mathbf{x} + \Delta \mathbf{x}) - F(\mathbf{x}) \\ \leq -\frac{1}{2} \sum_{j \in \mathcal{P}_t} (\delta x_j)^2 + \frac{1}{2} \sum_{\substack{i,j \in \mathcal{P}_t, \\ i \neq j}} (\mathbf{A}^T \mathbf{A})_{i,j} \delta x_i \delta x_j \end{aligned} \quad (35)$$

6. Detailed Proof: Theorem 3.2

This proof uses the result from Lemma 3.3, which is proven in detail in Sec. 7.

Modify the potential function used for sequential SCD:

$$\Psi(\mathbf{x}) = \frac{\beta}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 + \frac{1}{1-\epsilon} F(\mathbf{x}), \quad (36)$$

where ϵ is defined as in Eq. (67). Assume that \mathbf{P} is chosen s.t. $\epsilon < 1$.

Write out the change in the potential function from an update $\Delta \mathbf{x}$:

$$\begin{aligned} \Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x}) \\ = \frac{\beta}{2} [\|\mathbf{x} - \mathbf{x}^*\|_2^2 - \|\mathbf{x} + \Delta \mathbf{x} - \mathbf{x}^*\|_2^2] \\ + \frac{1}{1-\epsilon} [F(\mathbf{x}) - F(\mathbf{x} + \Delta \mathbf{x})] \end{aligned} \quad (37)$$

$$\begin{aligned} = \frac{\beta}{2} [-2\mathbf{x}^T \Delta \mathbf{x} + 2\mathbf{x}^{*T} \Delta \mathbf{x} - (\Delta \mathbf{x})^T (\Delta \mathbf{x})] \\ + \frac{1}{1-\epsilon} [F(\mathbf{x}) - F(\mathbf{x} + \Delta \mathbf{x})] \end{aligned} \quad (38)$$

$$\begin{aligned} = \beta \left[\sum_{j \in \mathcal{P}_t} -x_j \delta x_j + x_j^* \delta x_j - \frac{(\delta x_j)^2}{2} \right] \\ + \frac{1}{1-\epsilon} [F(\mathbf{x}) - F(\mathbf{x} + \Delta \mathbf{x})] \end{aligned} \quad (39)$$

Take the expectation w.r.t. \mathcal{P}_t , and use Lemma 3.3:

$$\begin{aligned} \mathbf{E}_{\mathcal{P}_t} [\Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x})] \\ = \beta \mathbf{P} \mathbf{E}_j \left[-x_j \delta x_j + x_j^* \delta x_j - \frac{(\delta x_j)^2}{2} \right] \\ + \frac{1}{1-\epsilon} \mathbf{E}_{\mathcal{P}_t} [F(\mathbf{x}) - F(\mathbf{x} + \Delta \mathbf{x})] \end{aligned} \quad (40)$$

$$\begin{aligned} \geq \beta \mathbf{P} \mathbf{E}_j \left[-x_j \delta x_j + x_j^* \delta x_j - \frac{(\delta x_j)^2}{2} \right] \\ - \mathbf{P} \frac{1}{1-\epsilon} \mathbf{E}_j \left[\delta x_j (\nabla F(\mathbf{x}))_j + \frac{\beta}{2} (1+\epsilon) (\delta x_j)^2 \right] \end{aligned} \quad (41)$$

$$\begin{aligned} = \beta \mathbf{P} \mathbf{E}_j \left[-x_j \delta x_j + x_j^* \delta x_j - \frac{1}{1-\epsilon} (\delta x_j)^2 \right] \\ - \mathbf{P} \frac{1}{1-\epsilon} \mathbf{E}_j [\delta x_j (\nabla F(\mathbf{x}))_j] \end{aligned} \quad (42)$$

$$\begin{aligned} \geq \beta \mathbf{P} \mathbf{E}_j \left[-x_j \delta x_j - x_j^* (\nabla F(\mathbf{x}))_j / \beta \right. \\ \left. - \frac{1}{1-\epsilon} (\delta x_j)^2 - \frac{1}{1-\epsilon} \delta x_j (\nabla F(\mathbf{x}))_j / \beta \right] \end{aligned} \quad (43)$$

where the last inequality used the update rule in Eq. (4), which implies $\delta x_j \geq -(\nabla F(\mathbf{x}))_j / \beta$.

Consider the two cases in the update rule in Eq. (4). Case 1: $\delta x_j = -x_j \geq -(\nabla F(\mathbf{x}))_j / \beta$.

$$\begin{aligned} \mathbf{E}_{\mathcal{P}_t} [\Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x})] \\ \geq \beta \mathbf{P} \mathbf{E}_j \left[-x_j \delta x_j - x_j^* (\nabla F(\mathbf{x}))_j / \beta \right. \\ \left. + \frac{1}{1-\epsilon} x_j \delta x_j + \frac{1}{1-\epsilon} x_j (\nabla F(\mathbf{x}))_j / \beta \right] \end{aligned} \quad (44)$$

$$\begin{aligned} = \beta \mathbf{P} \mathbf{E}_j \left[\frac{\epsilon}{1-\epsilon} x_j \delta x_j - x_j^* (\nabla F(\mathbf{x}))_j / \beta \right. \\ \left. + \frac{1}{1-\epsilon} x_j (\nabla F(\mathbf{x}))_j / \beta \right] \end{aligned} \quad (45)$$

$$\begin{aligned} \geq \beta \mathbf{P} \mathbf{E}_j \left[-\frac{\epsilon}{1-\epsilon} x_j (\nabla F(\mathbf{x}))_j / \beta - x_j^* (\nabla F(\mathbf{x}))_j / \beta \right. \\ \left. + \frac{1}{1-\epsilon} x_j (\nabla F(\mathbf{x}))_j / \beta \right] \end{aligned} \quad (46)$$

$$= \mathbf{P} \mathbf{E}_j [(x_j - x_j^*) (\nabla F(\mathbf{x}))_j]. \quad (47)$$

Case 2: $\delta x_j = -(\nabla F(\mathbf{x}))_j / \beta \geq -x_j$.

$$\begin{aligned} \mathbf{E}_{\mathcal{P}_t} [\Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x})] \\ \geq \beta \mathbf{P} \mathbf{E}_j \left[-x_j \delta x_j - x_j^* (\nabla F(\mathbf{x}))_j / \beta \right] \end{aligned} \quad (48)$$

$$\geq \mathbf{P} \mathbf{E}_j [(x_j - x_j^*) (\nabla F(\mathbf{x}))_j]. \quad (49)$$

In both cases,

$$\begin{aligned} \mathbf{E}_{\mathcal{P}_t} [\Psi(\mathbf{x}) - \Psi(\mathbf{x} + \Delta \mathbf{x})] \\ \geq \mathbf{P} \mathbf{E}_j [(x_j - x_j^*) (\nabla F(\mathbf{x}))_j] \end{aligned} \quad (50)$$

$$= \frac{\mathbf{P}}{2d} (\mathbf{x} - \mathbf{x}^*)^T \nabla F(\mathbf{x}) \quad (51)$$

$$\geq \frac{\mathbf{P}}{2d} (F(\mathbf{x}) - F(\mathbf{x}^*)), \quad (52)$$

where the last inequality holds since $F(\mathbf{x})$ is convex.

Now sum over $T + 1$ iterations (with an expectation over the \mathcal{P}_t from all iterations):

$$\mathbf{E} \left[\sum_{t=0}^T \Psi(\mathbf{x}^{(t)}) - \Psi(\mathbf{x}^{(t+1)}) \right] \geq \frac{\mathbb{P}}{2d} \mathbf{E} \left[\sum_{t=0}^T F(\mathbf{x}^{(t)}) - F(\mathbf{x}^*) \right] \quad (53)$$

$$= \frac{\mathbb{P}}{2d} \left[\mathbf{E} \left[\sum_{t=0}^T F(\mathbf{x}^{(t)}) \right] - (T+1)F(\mathbf{x}^*) \right] \quad (54)$$

$$\geq \frac{\mathbb{P}(T+1)}{2d} \left[\mathbf{E} \left[F(\mathbf{x}^{(T)}) \right] - F(\mathbf{x}^*) \right], \quad (55)$$

where Eq. (55) uses the result from Lemma 3.3, which implies that the objective is decreasing in expectation for \mathbb{P} s.t. $\epsilon \leq 1$. (To see why the objective decreases in expectation, plug in the update rule in Eq. (4) into Eq. (68), and note that the right-hand side of Eq. (68) is negative.)

Since $\sum_{t=0}^T \Psi(\mathbf{x}^{(t)}) - \Psi(\mathbf{x}^{(t+1)}) = \Psi(\mathbf{x}^{(0)}) - \Psi(\mathbf{x}^{(T+1)})$, rearranging the above inequality gives

$$\mathbf{E} \left[F(\mathbf{x}^{(T)}) \right] - F(\mathbf{x}^*) \leq \frac{2d}{\mathbb{P}(T+1)} \mathbf{E} \left[\Psi(\mathbf{x}^{(0)}) - \Psi(\mathbf{x}^{(T+1)}) \right] \quad (56)$$

$$\leq \frac{2d}{\mathbb{P}(T+1)} \mathbf{E} \left[\Psi(\mathbf{x}^{(0)}) \right] \quad (57)$$

$$= \frac{2d}{\mathbb{P}(T+1)} \left[\frac{\beta}{2} \|\mathbf{x}^*\|_2^2 + \frac{1}{1-\epsilon} F(\mathbf{x}^{(0)}) \right]. \quad (58)$$

7. Detailed Proof: Lemma 3.3

Note: Assume the algorithm chooses a set of \mathbb{P} coordinates, not a multiset.

Starting with Assumption 3.1, we can rearrange terms as follows:

$$F(\mathbf{x} + \Delta\mathbf{x}) - F(\mathbf{x}) \leq (\Delta\mathbf{x})^T \nabla F(\mathbf{x}) + \frac{\beta}{2} (\Delta\mathbf{x})^T \mathbf{A}^T \mathbf{A} (\Delta\mathbf{x}) \quad (59)$$

Take the expectation w.r.t. \mathcal{P}_t , the set of updated weights, and use the fact that each set \mathcal{P}_t is equally

likely to be chosen.

$$\mathbf{E}_{\mathcal{P}_t} [F(\mathbf{x} + \Delta\mathbf{x}) - F(\mathbf{x})] \leq \mathbf{E}_{\mathcal{P}_t} \left[\sum_{j \in \mathcal{P}_t} \delta x_j (\nabla F(\mathbf{x}))_j \right] \quad (60)$$

$$+ \frac{\beta}{2} \mathbf{E}_{\mathcal{P}_t} \left[\sum_{i,j \in \mathcal{P}_t} \delta x_i (\mathbf{A}^T \mathbf{A})_{i,j} \delta x_j \right] = \mathbf{E}_{\mathcal{P}_t} \left[\sum_{j \in \mathcal{P}_t} \delta x_j (\nabla F(\mathbf{x}))_j + \frac{\beta}{2} (\delta x_j)^2 \right] \quad (61)$$

$$+ \frac{\beta}{2} \mathbf{E}_{\mathcal{P}_t} \left[\sum_{\substack{i,j \in \mathcal{P}_t \\ i \neq j}} \delta x_i (\mathbf{A}^T \mathbf{A})_{i,j} \delta x_j \right] = \mathbb{P} \mathbf{E}_j \left[\delta x_j (\nabla F(\mathbf{x}))_j + \frac{\beta}{2} (\delta x_j)^2 \right] \quad (62)$$

$$+ \frac{\beta}{2} \mathbb{P}(\mathbb{P}-1) \mathbf{E}_{i,j:i \neq j} \left[\delta x_i (\mathbf{A}^T \mathbf{A})_{i,j} \delta x_j \right],$$

where $\mathbf{E}_j[\cdot]$ denotes an expectation w.r.t. j chosen uniformly at random from $\{1, \dots, 2d\}$ and where $\mathbf{E}_{i,j:i \neq j}[\cdot]$ denotes an expectation w.r.t. a pair of distinct indices i, j chosen uniformly at random from $\{1, \dots, 2d\}$.

Since indices in \mathcal{P}_t are chosen uniformly at random, the expectations may be rewritten as

$$\mathbf{E}_{\mathcal{P}_t} [F(\mathbf{x} + \Delta\mathbf{x}) - F(\mathbf{x})] \quad (63)$$

$$\leq \frac{\mathbb{P}}{2d} \left[(\Delta\mathbf{x})^T (\nabla F(\mathbf{x})) + \frac{\beta}{2} (\Delta\mathbf{x})^T (\Delta\mathbf{x}) \right] \quad (64)$$

$$+ \frac{\beta}{2} \frac{\mathbb{P}(\mathbb{P}-1)}{2d(2d-1)} \left[(\Delta\mathbf{x})^T \mathbf{A}^T \mathbf{A} (\Delta\mathbf{x}) - (\Delta\mathbf{x})^T (\Delta\mathbf{x}) \right],$$

where we are overloading the notation $\Delta\mathbf{x}$: in Eq. (63), $\Delta\mathbf{x}$ only has non-zero entries in elements indexed by \mathcal{P}_t ; in Eq. (64), $\Delta\mathbf{x}$ can have non-zero entries everywhere (set by the update rule in Eq. (4)).

The spectral radius, i.e., the largest eigenvalue, of a matrix \mathbf{M} may be expressed as $\max_{\mathbf{z}} \frac{\mathbf{z}^T \mathbf{M} \mathbf{z}}{\mathbf{z}^T \mathbf{z}}$; see, e.g., Strang (?). Letting ρ be the spectral radius of $\mathbf{A}^T \mathbf{A}$, upper-bound the second-order term in Eq. (64):

$$\mathbf{E}_{\mathcal{P}_t} [F(\mathbf{x} + \Delta\mathbf{x}) - F(\mathbf{x})] \leq \frac{\mathbb{P}}{2d} \left[(\Delta\mathbf{x})^T (\nabla F(\mathbf{x})) + \frac{\beta}{2} (\Delta\mathbf{x})^T (\Delta\mathbf{x}) \right] \quad (65)$$

$$+ \frac{\beta}{2} \frac{\mathbb{P}(\mathbb{P}-1)}{2d(2d-1)} \left[\rho (\Delta\mathbf{x})^T (\Delta\mathbf{x}) - (\Delta\mathbf{x})^T (\Delta\mathbf{x}) \right] = \frac{\mathbb{P}}{2d} (\Delta\mathbf{x})^T (\nabla F(\mathbf{x})) \quad (66)$$

$$+ \frac{\beta}{2} \frac{\mathbb{P}}{2d} \left(1 + \frac{(\mathbb{P}-1)(\rho-1)}{2d-1} \right) (\Delta\mathbf{x})^T (\Delta\mathbf{x})$$

Letting

$$\epsilon = \frac{(\mathbb{P}-1)(\rho-1)}{2d-1}, \quad (67)$$

we can rewrite the right-hand side in terms of expectations over $j \in \{1, \dots, 2d\}$ to get the lemma's result:

$$\begin{aligned} & \mathbf{E}_{\mathcal{P}_t} [F(\mathbf{x} + \Delta\mathbf{x}) - F(\mathbf{x})] \\ & \leq \mathbf{P}\mathbf{E}_j \left[\delta x_j (\nabla F(\mathbf{x}))_j + \frac{\beta}{2} (1 + \epsilon) (\delta x_j)^2 \right]. \end{aligned} \quad (68)$$

Note: If we let the algorithm choose a multiset, rather than a set, of \mathbf{P} coordinates, then we get $\epsilon = \frac{(\mathbf{P}-1)\rho}{2d}$, which gives worse scaling than the ϵ above. (Compare the two when all features are uncorrelated so that $\rho = 1$. With a set, the ϵ above indicates that we can use $\mathbf{P} = 2d$ and get good scaling; with a multiset, the changed ϵ indicates that we can be hurt by using larger \mathbf{P} .)