

7 Appendix

7.1 Additional Results for Section 3

The following lemma provides upper bounds on the expected gradient of the worst-possible MKL-SGD solution that lies in a ball around \mathbf{w}^* . Simultaneously satisfying the following bound with the one in Lemma 3 may lead to an infeasible set of ϵ and N' . And thus we use Lemma 4 in conjunction with 3.

Lemma 6. *Let us assume that MKL-SGD converges to $\bar{\mathbf{w}}_{MKL}$. For any $\bar{\mathbf{w}}_{MKL} \in \mathcal{B}_r(\mathbf{w}^*)$ that satisfies assumptions N1, N2, A4 and A5, there exists $N' \geq N$ and $\epsilon' \leq \epsilon$ such that,*

$$\left\| \sum_{i \notin \mathbb{O}} p_i(\bar{\mathbf{w}}_{MKL}) \nabla f_i(\bar{\mathbf{w}}_{MKL}) \right\| \leq \min \left\{ (1 - \epsilon^k) L \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\|, \epsilon^k G(\mathbf{w}) \right\}$$

The proof for lemma 2 can be found in the Appendix Section 7.2.8

7.1.1 Squared loss in the scalar setting with outliers centered at different points

We will assume that without loss of generality all the outliers will lie on the same side of w^* . If that's not the case, the bounds which show in the subsequent part will be even stronger. Without loss of generality, assume $0 < w_{b_1} < w_{b_2} < \dots < w_{b_{|\mathbb{B}|}}$.

The loss functions and \tilde{w} are redefined as follows:

$$f_i(w) = \begin{cases} l_i(w - w^*)^2 & \forall i \notin \mathbb{O} \\ l_i(w - w_{b_i})^2 & \forall i \in \mathbb{O}, \end{cases} \quad (14)$$

$$\tilde{w} := \left\{ w \mid w = \min_{\lambda \in (0,1)} \left[\lambda w^* + (1 - \lambda) \min_j w_{b_j} \right], f_{l_m}(w) = f_{l_M}(w) \right\} \quad (15)$$

Once again by simple analysis of different points of intersection, we can describe a closed form expression of \tilde{w} as follows:

$$\tilde{w} = \frac{\sqrt{l_m} w^* + \sqrt{l_{\tilde{M}}} w_{b_{\tilde{M}}}}{\sqrt{l_m} + \sqrt{l_{\tilde{M}}}}$$

where $\tilde{M} = \arg \max_{j \in \mathbb{O}} \sqrt{l_j} w_{b_j}$ and $\kappa = \frac{l_M}{l_m}$ and $\rho = \frac{\min_i w_{b_i}}{\max_i w_{b_i}}$

Condition 2. $\hat{p} < \frac{1}{1 + \frac{1}{\rho \left(1 + \frac{1}{\sqrt{\kappa}} \right)} - 1}$

When $\rho = 1$, condition 2 becomes identical to condition 1.

Lemma 7. *If Condition 2 is satisfied and the loss functions and \tilde{w} are defined as in equation 26. Now let us start at a point where the highest probabilities are assigned to all the bad samples, even in that case, the stationary point attained by MKL-SGD will be such that the highest probabilities are assigned to the good samples.*

7.2 Proofs and supporting lemmas

7.2.1 Proof of Lemma 1

Proof. Proof. $\tilde{F}(\mathbf{w}) = \sum_i p_{m_i(\mathbf{w})}(\mathbf{w})$. Let us fix a \mathbf{w} such that $p_i = p_i(\mathbf{w})$. We know that for any p_i , $\sum_i p_i f_i(\mathbf{w})$ is strongly convex in \mathbf{w} with parameter $\lambda_{\mathbf{w}}$. This implies

$$\nabla \tilde{F}(\mathbf{w})^\top (\mathbf{w} - \mathbf{w}^*) \geq \lambda_{\mathbf{w}} \|\mathbf{w} - \mathbf{w}^*\|^2$$

□

A naive bound for the above Lemma can be:

$$\nabla \tilde{F}(\mathbf{w})^\top (\mathbf{w} - \mathbf{w}^*) \geq \min_i p_i \sum_i f_i(\mathbf{w}) \geq \underbrace{\lambda \min_i p_i}_{\lambda_{\mathbf{w}}} \|\mathbf{w} - \mathbf{w}^*\|^2$$

□

7.2.2 Proof of Theorem 1

Proof. By the definition of the noiseless framework, \mathbf{w}^* is the unique optimum of $F(\mathbf{w})$ and lies in the optimal set of each $f_i(\cdot)$. We will prove this theorem by contradiction. Assume there exists some $\hat{\mathbf{w}} \neq \mathbf{w}^*$ that also satisfies optimum of $\nabla \tilde{F}(\hat{\mathbf{w}}) = \mathbf{0}$. At $\hat{\mathbf{w}}$, we have $0 = \langle \nabla \tilde{F}(\hat{\mathbf{w}}), \hat{\mathbf{w}} - \mathbf{w}^* \rangle = \lambda \|\hat{\mathbf{w}} - \mathbf{w}^*\|^2$. This implies $\hat{\mathbf{w}} = \mathbf{w}^*$. □

Theorem 1 and Assumption 2 guarantee that $\lambda_{\mathbf{w}} > 0$. If $f(\mathbf{w})$ is strongly convex and $g(\mathbf{w})$ is convex, then we know that $f(\mathbf{w}) + g(\mathbf{w})$ is strongly convex. On similar lines we can show that $\lambda > 0$ by splitting the terms in $\tilde{F}(\mathbf{w})$ as $p_{\min} F(\mathbf{w})$ and $(\tilde{F}(\mathbf{w}) - p_{\min} F(\mathbf{w}))$. The first term has $\lambda > 0$ (Assumption 2) and the second term has $\lambda = 0$ (since it is convex). Note, p_{\min} is a positive constant independent of \mathbf{w} and so the above lemma is for all \mathbf{w} .

7.2.3 Proof of the claim in Section 3

We will first describe the problem setting again for ease of analysis before elaborating on the proof.

Problem setting Let us assume good and bad samples are centered at the same point with different Lipschitz constants. The loss functions are given as follows:

$$f_i(w) = \begin{cases} l_i(w - w^*)^2 & \forall i \notin \mathbb{O} \\ l_i(w - w_B)^2 & \forall i \in \mathbb{O}, \end{cases} \quad (16)$$

where $|\mathbb{O}| = b$ such that $n = g + b$. Let $l_m = \min_{i \notin \mathbb{O}} l_i$ and Let $l_M = \max_{i \in \mathbb{O}} l_i$ and $l_{max} = \min_{i \in [n]} l_i$, $l_{min} = \min_{i \in [n]} l_i$. Let us define $\kappa = \frac{l_{max}}{l_{min}} \geq \frac{l_M}{l_m}$. Let us define \tilde{w} as follows:

$$\tilde{w} := \left\{ w \mid \begin{array}{l} w = \min_{\alpha} \alpha w^* + (1 - \alpha) w_B, \quad \alpha \in (0, 1), \\ f_{l_m}(w) = f_{l_M}(w) \end{array} \right\} \quad (17)$$

By observation, we know for the scalar case $\tilde{w} = \frac{\sqrt{l_m} w^* + \sqrt{l_M} w_B}{\sqrt{l_m} + \sqrt{l_M}}$. Since we are initializing at w_0 , the probability of picking bad samples is \hat{p} .

Proof. Assume we start at w_B , such that the outlier functions will have the highest weights. Without loss of generality, the probability of picking sample j in the bad set is $p_{m_j(w_B)}(w_B)$.

Let \bar{w} indicate the stationary point of MKL-SGD assuming fixed $p_i(w_B)$ centered at w_B . The probabilities will not change until we reach \bar{w} . At \bar{w} , we have:

$$\begin{aligned} \sum_{i \notin \mathbb{O}} p_i(w_B) \nabla f_i(\bar{w}) &= - \sum_{j \in \mathbb{O}} p_j(w_B) \nabla f_j(\bar{w}) \\ \sum_{i \notin \mathbb{O}} p_i(w_B) l_i(\bar{w} - w^*) &= - \sum_{j \in \mathbb{O}} p_j(w_B) l_j(\bar{w} - w_B) \\ \bar{w} &= \frac{\sum_{i \notin \mathbb{O}} p_i(w_B) l_i w^* + \sum_{j \in \mathbb{O}} p_j(w_B) l_j w_B}{\sum_{i \notin \mathbb{O}} p_i(w_B) l_i + \sum_{j \in \mathbb{O}} p_j(w_B) l_j} \end{aligned}$$

If \bar{w} is closer to w_B than \tilde{w} , then the local minima \bar{w} exists, else we will escape that local minima, since at \tilde{w} , the probabilities change.

When is $\bar{w} < \tilde{w}$?

A sufficient condition for that is:

$$\begin{aligned} \frac{\sum_{j \in \mathbb{O}} p_j(w_B) l_j w_B}{\sum_{i \notin \mathbb{O}} p_i(w_B) l_i + \sum_{j \in \mathbb{O}} p_j(w_B) l_j} &\leq \frac{\sqrt{l_M} w_B}{\sqrt{l_m} + \sqrt{l_M}} \\ \frac{1}{1 + \frac{\sum_{i \notin \mathbb{O}} p_i(w_B) l_i}{\sum_{j \in \mathbb{O}} p_j(w_B) l_j}} &\leq \frac{1}{1 + \frac{\sqrt{l_m}}{\sqrt{l_M}}} \\ \frac{\sqrt{l_m}}{\sqrt{l_M}} &\leq \frac{\sum_{i \notin \mathbb{O}} p_i(w_B) l_i}{\sum_{j \in \mathbb{O}} p_j(w_B) l_j} \end{aligned}$$

We just need,

$$\begin{aligned} \frac{\sqrt{l_m}}{\sqrt{l_M}} &\leq \frac{\sum_{i \notin \mathbb{O}} p_i(w_B) l_m}{\sum_{j \in \mathbb{O}} p_j(w_B) l_M} \\ \frac{\sqrt{l_m}}{\sqrt{l_M}} &\leq \frac{(1 - \hat{p}) l_m}{\hat{p} l_M} \\ \frac{\hat{p}}{1 - \hat{p}} &\leq \sqrt{\frac{l_m}{l_M}} \\ \hat{p} &\leq \frac{1}{1 + \sqrt{\frac{l_M}{l_m}}} \end{aligned}$$

If $\hat{p} \leq \frac{1}{1 + \sqrt{\kappa}}$, then the above inequality is satisfied.

This was the condition for the point of intersection between the curves of good sample with the smallest l and the bad sample with the largest l . What happens between w^* and \tilde{w} ?

Next, we evaluate the closest point of intersection to \tilde{w} between w^* and \tilde{w} and do that recursively.

$$\begin{aligned} \frac{\sqrt{l_{m1}}}{\sqrt{l_{M1}}} &\leq \frac{(1 - \hat{p}) l_m}{\hat{p} l_M} \\ \hat{p} &\leq \frac{1}{1 + \sqrt{\frac{l_M}{l_m}} \sqrt{\frac{l_M}{l_{M1}}} \sqrt{\frac{l_{m1}}{l_m}}} \end{aligned}$$

Similarly, if $\hat{p} \leq \frac{1}{1 + \sqrt{\kappa}^{1.5}}$, then the above inequality is satisfied. □

7.2.4 Proof of Lemma 2

Let \bar{w} be a stationary point of MKL-SGD. Now, we analyze the loss landscape on the line joining w^* and w_C where $w_C = C\bar{w}$ is any arbitrary point ⁴ in the landscape at a distance as far as the farthest outlier from w^* . Let C be a very large number.

⁴Note that we just need w_C for the purpose of landscape analysis and it is not a parameter of the algorithm

The loss functions and $\tilde{\mathbf{w}}$ are redefined as follows:

$$f_i(\mathbf{w}) = \begin{cases} l_i \|\mathbf{w} - \mathbf{w}^*\|^2 & \forall i \in \mathbb{O} \\ l_i \|\mathbf{w} - \mathbf{w}_{b_i}\|^2 & \forall i \notin \mathbb{O}, \end{cases}$$

$$\tilde{\mathbf{w}} := \left\{ \mathbf{w} \mid \begin{array}{l} \mathbf{w} = \min_{\alpha \in (0,1)} \alpha \mathbf{w}^* + (1-\alpha) \mathbf{w}_C, \\ f_{l_m}(\mathbf{w}) = f_{l_M}(\mathbf{w}) \end{array} \right\}$$

where $|\mathbb{O}| = b$ such that $n = g + b$. Let $l_m = \min_{i \notin \mathbb{O}} l_i$ and Let $l_M = \max_{i \in \mathbb{O}} l_i$ and $l_{max} = \min_{i \in [n]} l_i$, $l_{min} = \min_{i \in [n]} l_i$. Let us define $\kappa = \frac{l_{max}}{l_{min}} \geq \frac{l_M}{l_m}$.

Now at $\bar{\mathbf{w}}$, we have $\nabla \tilde{F}(\bar{\mathbf{w}}) = 0$. Let us assume that the outliers are chosen in such a way that at \mathbf{w}_C , all the outliers have the lowest loss. As stated in the previous lemma, the results hold irrespective of that. This implies:

$$\begin{aligned} \sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) \nabla f_i(\bar{\mathbf{w}}) &= - \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) \nabla f_j(\bar{\mathbf{w}}) \\ \sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i (\bar{\mathbf{w}} - \mathbf{w}^*) &= - \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j (\bar{\mathbf{w}} - \mathbf{w}_{b_j}) \\ \bar{\mathbf{w}} &= \frac{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i \mathbf{w}^* + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j \mathbf{w}_{b_j}}{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \\ \text{By triangle inequality, } \|\bar{\mathbf{w}} - \mathbf{w}^*\| &\leq \frac{\sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j \|\mathbf{w}_{b_j} - \mathbf{w}^*\|}{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \end{aligned}$$

Without loss of generality assume that the outliers are ordered as follows: $\|\mathbf{w}_{b_1} - \mathbf{w}^*\| \leq \|\mathbf{w}_{b_2} - \mathbf{w}^*\| \leq \dots \leq \|\mathbf{w}_{b_{|\mathbb{O}|}} - \mathbf{w}^*\|$.

Now $\tilde{\mathbf{w}}$ be some point of intersection of function in the set of clean samples and a function in the set of outliers to \mathbf{w}^* . Let θ_j be the angle between the line connecting \mathbf{w}_{b_j} and \mathbf{w}^* to the line connecting \mathbf{w}_C to \mathbf{w}^* . For any two curves with Lipschitz constants l_i and l_j , the halfspaces passing through the weighted mean are also the region where both functions have equal values.

Thus,

$$\tilde{\mathbf{w}} = \frac{\sqrt{l_i} \mathbf{w}^* + \sqrt{l_j} \mathbf{w}_{b_j}}{\sqrt{l_i} + \sqrt{l_j}}$$

.

$$\|\tilde{\mathbf{w}} - \mathbf{w}^*\| = \frac{\sqrt{l_j} \|\mathbf{w}_{b_j} - \mathbf{w}^*\|}{\sqrt{l_j} + \sqrt{l_i}}$$

Let γ denote the following ratio:

$$\gamma = \frac{\min_{j \in \mathbb{O}} \|\mathbf{w}_{b_j} - \mathbf{w}^*\|}{\max_{j \in \mathbb{O}} \|\mathbf{w}_{b_j} - \mathbf{w}^*\|} = \frac{2\delta}{\delta_{max}}$$

Now, we want:

$$\begin{aligned}
 & \frac{\sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j \|\mathbf{w}_{b_j} - \mathbf{w}^*\|}{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \leq \frac{\sqrt{l_{t_j}}}{\sqrt{l_{t_j}} + \sqrt{l_g}} \frac{\|\mathbf{w}_{b_j} - \mathbf{w}^*\|}{\cos \theta_j} = \frac{\|\tilde{\mathbf{w}} - \mathbf{w}^*\|}{\cos \theta_j} \\
 & \frac{\sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j \|\mathbf{w}_{b_j} - \mathbf{w}^*\|}{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \leq \frac{\sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j \|\mathbf{w}_{b_{|O|}} - \mathbf{w}^*\|}{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \leq \frac{\sqrt{l_{t_j}}}{\sqrt{l_{t_j}} + \sqrt{l_g}} \frac{\|\mathbf{w}_{b_j} - \mathbf{w}^*\|}{\cos \theta_j} \\
 & \frac{\sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j}{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \leq \frac{\sqrt{l_{t_j}}}{\sqrt{l_{t_j}} + \sqrt{l_g}} \frac{\|\mathbf{w}_{b_j} - \mathbf{w}^*\|}{\cos \theta_j \|\mathbf{w}_{b_{|O|}} - \mathbf{w}^*\|} \\
 & \frac{\sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j}{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \leq \frac{\sqrt{l_{t_j}}}{\sqrt{l_{t_j}} + \sqrt{l_g}} \frac{\gamma}{\cos \theta_j}
 \end{aligned}$$

For simplicity, $\Gamma = \frac{\gamma}{\cos \theta_j}$, then we have:

$$\begin{aligned}
 & \frac{\sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j}{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i + \sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \leq \frac{\sqrt{l_{t_j}}}{\sqrt{l_{t_j}} + \sqrt{l_g}} \Gamma \\
 & \frac{1}{\Gamma} \left(\frac{\sqrt{l_g}}{\sqrt{l_{t_j}}} + 1 \right) - 1 \leq \frac{(1 - \hat{p}) l_m}{\hat{p} l_M} \leq \frac{\sum_{i \notin \mathbb{O}} p_i(\mathbf{w}_C) l_i}{\sum_{j \in \mathbb{O}} p_j(\mathbf{w}_C) l_j} \\
 & \frac{\hat{p}}{1 - \hat{p}} \leq \frac{\frac{l_m}{l_M}}{\frac{1}{\Gamma} - 1 + \frac{1}{\Gamma} \frac{\sqrt{l_g}}{\sqrt{l_{t_j}}}} \\
 & \hat{p} \leq \frac{1}{1 + \kappa \left(\frac{1}{\Gamma} - 1 + \frac{\sqrt{\kappa}}{\Gamma} \right)} \leq \frac{1}{1 + \frac{l_M}{l_m} \left(\frac{1}{\Gamma} - 1 + \frac{1}{\Gamma} \frac{\sqrt{l_g}}{\sqrt{l_{t_j}}} \right)}
 \end{aligned}$$

Replacing $\Gamma = \frac{\gamma}{\cos \theta_j}$, and let $q = \frac{\cos \theta_j}{\gamma} - 1 + \frac{\cos \theta_j \sqrt{\kappa}}{\gamma}$ the condition to guarantee that bad local minima do not exist is $\hat{p} \leq \frac{1}{1 + \kappa q}$ and $q > 0$.

Note: In the vector case, for example there exists a fine tradeoff between how large θ_j can be and if for large θ_j , the loss corresponding to the outlier will be one of the lowest. Understanding that tradeoff is beyond the scope of this paper.

Note that, the lemma 2 leads to a very strong worst-case guarantee. It states that the farthest optimum will always be within a bowl of distance r from \mathbf{w}^* no matter where we initialize. Moreover, as long as the condition is satisfied no matter where the outliers lie (can be adversarially chosen), MKL-SGD always has the propensity to bring the iterates to a ball of radius r around \mathbf{w}^* . However, when the necessary conditions for its convergence are violated, the guarantees are initialization dependent. Thus, all the discussions in the rest of this section will be with respect to these worst case guarantees. However, as we see in the experimental section for both neural networks and linear regression, random initialization also seems to perform better than SGD.

Effect of κ A direct result of Lemma 2 is that higher the condition number of the set of quadratic loss functions, lower is the fraction ϵ of outliers the MKL-SGD can tolerate. This is because large κ results in a small value of $\frac{1}{1 + \kappa q}$. This implies that \hat{p} has to be small which in turn requires smaller fractions of corruptions, ϵ .

Effect of γ : The relative distance of the outliers from \mathbf{w}^* plays a critical role in the condition for Lemma 2. We know that $\gamma \in (0, 1]$. $\gamma = 1$ implies the outliers are equidistant from the optimum \mathbf{w}^* . Low values

of γ lead to a large q leading to the violation of the condition with \hat{p} (since RHS in the condition is very small), which implies that one bad outlier can guarantee that the condition in Lemma 2 are violated. The guarantees in the above lemma are only when the outliers are not adversarially chosen to lie at very high relative distances from \mathbf{w}^* . One way to avoid the set of outliers far far away from the optimum is to have a filtering step at the start of the algorithm like the one in Diakonikolas et al. [2018]. We will refer this in Experiments.

Effect of $\cos \theta_{j, \bar{\mathbf{w}}}$: At first glance, it may seem that $\cos \theta_{j, \bar{\mathbf{w}}} = 0$ may cause $1 + \kappa q < 0$ and since $\hat{p}(\mathbf{w}) > 0$, the condition in Lemma 2 may never be satisfied. Since, the term $\cos \theta_{j, \bar{\mathbf{w}}}$ shows up in the denominator of the loss associated with outlier centered at \mathbf{w}_{b_j} . Thus, low values of $\cos \theta_{j, \bar{\mathbf{w}}}$ implies high value of loss associated with the function centered at \mathbf{w}_{b_j} which in turn implies the maximum probability attained by that sample can never be in the top- $|\mathbb{O}|$ probabilities for that $\bar{\mathbf{w}}$.

7.2.5 Proof of Lemma 3

Proof. At $\bar{\mathbf{w}}_{SGD}$, $\nabla \tilde{F}(\bar{\mathbf{w}}_{SGD}) = 0$. Then,

$$\begin{aligned}
 \sum_{i \notin \mathbb{O}} \nabla f_i(\bar{\mathbf{w}}_{SGD}) &= - \sum_{i \in \mathbb{O}} \nabla f_i(\bar{\mathbf{w}}_{SGD}) \\
 \left\| \sum_{i \notin \mathbb{O}} \nabla f_i(\bar{\mathbf{w}}_{SGD}) \right\| &= \left\| \sum_{i \in \mathbb{O}} \nabla f_i(\bar{\mathbf{w}}_{SGD}) \right\| \\
 \left\| \sum_{i \notin \mathbb{O}} \nabla f_i(\bar{\mathbf{w}}_{SGD}) \right\| &\leq \sum_i \|\nabla f_i(\bar{\mathbf{w}}_{SGD})\| \\
 &\leq \sum_i L \|\bar{\mathbf{w}}_{SGD} - \mathbf{w}^*\| \\
 &= (1 - \epsilon)nL \|\bar{\mathbf{w}}_{SGD} - \mathbf{w}^*\| \tag{18}
 \end{aligned}$$

$$\begin{aligned}
 \left\| \sum_{i \in \mathbb{O}} \nabla f_i(\bar{\mathbf{w}}_{SGD}) \right\| &\leq \sum_{i \in \mathbb{O}} \|\nabla f_i(\bar{\mathbf{w}}_{SGD})\| \\
 &\leq \sum_{i \in \mathbb{O}} G(\bar{\mathbf{w}}_{SGD}) \\
 &\leq \epsilon G(\bar{\mathbf{w}}_{SGD}) \tag{19}
 \end{aligned}$$

$$\left\| \sum_{i \in \mathbb{O}} \nabla f_i(\bar{\mathbf{w}}_{SGD}) \right\| = \min(\epsilon n G(\bar{\mathbf{w}}_{SGD}), (1 - \epsilon)nL \|\bar{\mathbf{w}}_{SGD} - \mathbf{w}^*\|)$$

□

7.2.6 Proof of Lemma 4

Proof. At $\bar{\mathbf{w}}_{MKL}$, $\nabla \tilde{F}(\bar{\mathbf{w}}_{MKL}) = 0$. This implies

$$\sum_{i \notin \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \nabla f_i(\bar{\mathbf{w}}_{MKL}) = - \sum_{i \in \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \nabla f_i(\bar{\mathbf{w}}_{MKL})$$

Multiplying both sides by $(\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*)$

$$\begin{aligned} \sum_{i \notin \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \langle \nabla f_i(\bar{\mathbf{w}}_{MKL}), \bar{\mathbf{w}}_{MKL} - \mathbf{w}^* \rangle &= - \sum_{i \in \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \langle \nabla f_i(\bar{\mathbf{w}}_{MKL}), \bar{\mathbf{w}}_{MKL} - \mathbf{w}^* \rangle \quad (20) \\ \langle \nabla \tilde{F}_G(\bar{\mathbf{w}}_{MKL}), \bar{\mathbf{w}}_{MKL} - \mathbf{w}^* \rangle &= - \sum_{i \in \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \langle \nabla f_i(\bar{\mathbf{w}}_{MKL}), \bar{\mathbf{w}}_{MKL} - \mathbf{w}^* \rangle \end{aligned}$$

Lower bounding the LHS using Lemma 1 and $m = m(\bar{\mathbf{w}}_{MKL})^5$,

$$m \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\|^2 \leq \left| \langle \nabla \tilde{F}_G(\bar{\mathbf{w}}_{MKL}), \bar{\mathbf{w}}_{MKL} - \mathbf{w}^* \rangle \right| = LHS \quad (21)$$

$$RHS \leq \left| - \sum_{i \in \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \langle \nabla f_i(\bar{\mathbf{w}}_{MKL}), \bar{\mathbf{w}}_{MKL} - \mathbf{w}^* \rangle \right|$$

$$m \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\|^2 \leq \sum_{i \in \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \|\langle \nabla f_i(\bar{\mathbf{w}}_{MKL}), \bar{\mathbf{w}}_{MKL} - \mathbf{w}^* \rangle\|$$

$$m \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\|^2 \leq \sum_{i \in \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \|\nabla f_i(\bar{\mathbf{w}}_{MKL})\| \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\|$$

$$m \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\|^2 \leq \sum_{i \in \mathcal{O}} p_i(\bar{\mathbf{w}}_{MKL}) \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\| G(\bar{\mathbf{w}}_{SGD})$$

$$m \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\| \leq \epsilon^k G(\bar{\mathbf{w}}_{SGD}) \quad (23)$$

□

7.2.7 Proof of Theorem 2

Proof. There exists an $\epsilon' \leq \epsilon$ such that in Lemma 3, we have

$$(1 - \epsilon)L \|\bar{\mathbf{w}}_{SGD} - \mathbf{w}^*\| \geq \epsilon G(\bar{\mathbf{w}}_{SGD})$$

Combining above equation with Lemma 4, we get

$$\begin{aligned} (1 - \epsilon)L \|\bar{\mathbf{w}}_{SGD} - \mathbf{w}^*\| &\geq \epsilon G(\bar{\mathbf{w}}_{SGD}) \geq \epsilon \frac{\lambda}{\epsilon^2} \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\| \\ \Rightarrow \|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\| &\leq \frac{(1 - \epsilon)L \epsilon^{k-1}}{\lambda} \|\bar{\mathbf{w}}_{SGD} - \mathbf{w}^*\| \end{aligned}$$

Picking a large enough k , we can guarantee that $\frac{(1 - \epsilon)L \epsilon^{k-1}}{\lambda} < 1$ □

7.2.8 Proof of Lemma 6

Proof. From the definition of good samples in the noiseless setting, we know that $f_i(\mathbf{w}^*) = 0 \forall i \notin \mathcal{O}$. Similarly, for samples belonging to the outlier set, $f_i(\mathbf{w}^*) > 0 \forall i \in \mathcal{O}$. There exists a ball around the optimum of radius r such that $f_i(\mathbf{w}) \leq f_j(\mathbf{w}) \forall i \notin \mathcal{O}, j \in \mathcal{O}, \mathbf{w} \in \mathbb{O}_r(\mathbf{w}^*)$. Assume that $N' \geq N$ and $\epsilon' \leq \epsilon$, such that $\|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\| \leq r$.

At \bar{w}_{MKL} , $\nabla \tilde{F}(\bar{w}_{MKL}) = 0$. This implies

$$\begin{aligned} \sum_{i \notin \mathcal{O}} p_i(\bar{w}_{MKL}) \nabla f_i(\bar{w}_{MKL}) &= - \sum_{i \in \mathcal{O}} p_i(\bar{w}_{MKL}) \nabla f_i(\bar{w}_{MKL}) \\ \left\| \sum_{i \notin \mathcal{O}} p_i(\bar{w}_{MKL}) \nabla f_i(\bar{w}_{MKL}) \right\| &= \left\| \sum_{i \in \mathcal{O}} p_i(\bar{w}_{MKL}) \nabla f_i(\bar{w}_{MKL}) \right\| \\ \left\| \sum_{i \notin \mathcal{O}} p_i(\bar{w}_{MKL}) \nabla f_i(\bar{w}_{MKL}) \right\| &\leq \sum_i p_i(\bar{w}_{MKL}) \|\nabla f_i(\bar{w}_{MKL})\| \\ &\leq \sum_i p_i(\bar{w}_{MKL}) L \|\bar{w}_{MKL} - w^*\| \\ &= (1 - \epsilon^k) L \|\bar{w}_{MKL} - w^*\| \end{aligned} \quad (24)$$

$$\begin{aligned} \left\| \sum_{i \in \mathcal{O}} p_i(\bar{w}_{MKL}) \nabla f_i(\bar{w}_{MKL}) \right\| &\leq \sum_{i \in \mathcal{O}} p_i(\bar{w}_{MKL}) \|\nabla f_i(\bar{w}_{MKL})\| \\ &\leq \sum_{i \in \mathcal{O}} p_i(\bar{w}_{MKL}) G(\bar{w}_{MKL}) \\ &\leq \epsilon^k G(\bar{w}_{MKL}) \end{aligned} \quad (25)$$

□

7.2.9 Proof of Lemma 7

Problem Setting: We will assume that without loss of generality all the outliers will lie on the same side of w^* . If that's not the case, the bounds which show in the subsequent part will be even stronger. Without loss of generality, assume $0 < w_{b_1} < w_{b_2} < \dots < w_{b_{|B|}}$.

The loss functions and \tilde{w} are redefined as follows:

$$f_i(w) = \begin{cases} l_i(w - w^*)^2 & \forall i \notin \mathcal{O} \\ l_i(w - w_{b_i})^2 & \forall i \in \mathcal{O}, \end{cases} \quad (26)$$

$$\tilde{w} := \left\{ w \mid w = \min_{\lambda \in (0,1)} \left[\lambda w^* + (1 - \lambda) \min_j w_{b_j} \right], f_{l_m}(w) = f_{l_M}(w) \right\} \quad (27)$$

Once again by simple analysis of different points of intersection, we can describe a closed form expression of \tilde{w} as follows:

$$\tilde{w} = \frac{\sqrt{l_m} w^* + \sqrt{l_M} w_{b_M}}{\sqrt{l_m} + \sqrt{l_M}}$$

where $\tilde{M} = \arg \max_{j \in \mathcal{O}} \sqrt{l_j} w_{b_j}$, $M = \arg \max_{i \in \mathcal{O}} l_i$, $m = \arg \max_{i \notin \mathcal{O}} l_i$.

Here, $|\mathcal{G}| = g$ and $|\mathcal{O}| = b$ such that $n = g + b$. Let $l_m = \min_{i \notin \mathcal{O}} l_i$ and Let $l_M = \max_{i \in \mathcal{O}} l_i$ and $l_{max} = \min_{i \in [n]} l_i$, $l_{min} = \min_{i \in [n]} l_i$. Let us define $\kappa = \frac{l_{max}}{l_{min}} \geq \frac{l_M}{l_m}$. Let $\gamma = \frac{w_{b_1}}{w_{b_{|B|}}}$

Proof. Assume we start at w_{b_1} , such that the outlier functions will have the highest weights. If any one of the outlier functions does not have the top $|B|$ probabilities at w_{b_1} , then \hat{p} in the subsequent bounds will be smaller and will still satisfy the final condition. Without loss of generality, the probability of picking sample j in the bad set is $p_{m_j(w_{b_1})}(w_{b_1})$. Note that, we assume that \tilde{w} lies between w^* and w_{b_1} and the outlier functions have the highest weights at w_{b_1} . If that's not the case, our bounds are still satisfied since there will be atleast one term which will not be a part of \hat{p} and so \bar{w} will be even closer to w^* .

Let \bar{w} indicate the stationary point of MKL-SGD assuming fixed $p_i(w_B)$ centered at w_B . The probabilities will not change until we reach \tilde{w} . At \bar{w} , we have:

$$\begin{aligned} \sum_{i \notin \mathbb{O}} \nabla f_i(\bar{w}) &= - \sum_{j \in \mathbb{O}} \nabla f_j(\bar{w}) \\ \sum_{i \notin \mathbb{O}} p_i(w_{b_1}) l_i (\bar{w} - w^*) &= - \sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j (\bar{w} - w_{b_1}) \\ \bar{w} &= \frac{\sum_{i \notin \mathbb{O}} p_i(w_{b_1}) l_i w^* + \sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j w_{b_1}}{\sum_{i \notin \mathbb{O}} p_i(w_{b_1}) l_i + \sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j} \end{aligned}$$

If \bar{w} is closer to w_{b_1} than \tilde{w} , the local minima \bar{w} exists, since the probabilities don't change before we reach \tilde{w} starting from w_{b_1} moving towards w^* .

When is $\bar{w} < \tilde{w}$

A sufficient condition for that is:

$$\begin{aligned} \frac{\sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j w_{b_j}}{\sum_{i \notin \mathbb{O}} p_i(w_{b_1}) l_i + \sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j} &\leq \frac{\sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j w_{b_M}}{\sum_{i \notin \mathbb{O}} p_i(w_{b_1}) l_i + \sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j} \leq \frac{\sqrt{l_M} w_{b_M}}{\sqrt{l_m} + \sqrt{l_M}} \\ &\leq \frac{1}{1 + \frac{\sum_{i \notin \mathbb{O}} p_i(w_{b_1}) l_i}{\sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j}} \leq \frac{\gamma}{1 + \frac{\sqrt{l_m}}{\sqrt{l_M}}} \leq \frac{\frac{w_{b_M}}{w_{b_1}}}{1 + \frac{\sqrt{l_m}}{\sqrt{l_M}}} \\ &\leq \frac{\sqrt{l_m}}{\sqrt{l_M}} \leq \frac{\sum_{i \notin \mathbb{O}} p_i(w_{b_1}) l_i}{\sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j} \end{aligned}$$

We have,

$$\begin{aligned} 1 + \frac{\sqrt{l_m}}{\sqrt{l_M}} &\leq \gamma \left(1 + \frac{(1 - \hat{p}) l_m}{\hat{p} l_M} \right) \leq \gamma \left(1 + \frac{\sum_{i \notin \mathbb{O}} p_i(w_{b_1}) l_i}{\sum_{j \in \mathbb{O}} p_j(w_{b_1}) l_j} \right) \\ \frac{1}{\gamma} \left(1 + \frac{\sqrt{l_m}}{\sqrt{l_M}} \right) - 1 &\leq \frac{(1 - \hat{p}) l_m}{\hat{p} l_M} \\ \frac{\hat{p} l_M}{(1 - \hat{p}) l_m} &\leq \frac{1}{\frac{1}{\gamma} \left(1 + \frac{\sqrt{l_m}}{\sqrt{l_M}} \right) - 1} \end{aligned}$$

$$\begin{aligned} \frac{\hat{p}}{(1 - \hat{p})} &\leq \frac{l_m}{l_M} \frac{1}{\frac{1}{\gamma} \left(1 + \frac{\sqrt{l_m}}{\sqrt{l_M}} \right) - 1} \\ \hat{p} &\leq \frac{1}{l_M \left(\frac{1}{\gamma} \left(1 + \frac{\sqrt{l_m}}{\sqrt{l_M}} \right) - 1 \right) + \frac{1}{l_m}} \\ &\leq \frac{1}{1 + \frac{l_M}{l_m} \left(\frac{1}{\gamma} - 1 + \frac{\sqrt{l_m}}{\gamma \sqrt{l_M}} \right)} \end{aligned}$$

Thus, for the minima to not exist, we need: $\hat{p} \leq \frac{1}{1 + \kappa \left(\left(\frac{1}{\gamma} - 1 \right) + \frac{\sqrt{\kappa}}{\gamma} \right)}$

The condition for any other pair of l_i, l_j is similar. This is because if we replace m by m_1 in the square root,

the final step of the analysis simplifies to:

$$\hat{p} \leq \frac{1}{1 + \frac{l_M}{l_m} \left(\left(\frac{1}{\gamma} - 1 \right) + \sqrt{\frac{l_{m_1}}{l_M} \frac{1}{\gamma}} \right)} \quad (28)$$

If $\hat{p} \leq \frac{1}{1 + \kappa \left(\left(\frac{1}{\gamma} - 1 \right) + \frac{\sqrt{\kappa}}{\gamma} \right)}$, then the above equation is satisfied and thus the equivalent condition still remains the same. \square

7.3 Additional results and proofs for Section 5

Consider the sample size n with bad set(outlier) \mathbb{O} and good set \mathcal{G} such that $|\mathcal{G}| = n - |\mathbb{O}|$. Define

$$F_{good}(\mathbf{w}) = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} f_i(\mathbf{w}).$$

We assume:

(1) (Stationary Point) Assume \mathbf{w}^* is the solution for the average loss function of good sample such that

$$\nabla F_{good}(\mathbf{w}^*) = 0 \quad \text{but } \nabla f_i(\mathbf{w}^*) \neq 0, \forall i \in \mathbb{O}$$

(2) (Strong Convexity) $F_{good}(\mathbf{w})$ is strongly convex with parameters λ_{good} i.e.,

$$\langle \nabla F_{good}(\mathbf{w}) - \nabla F_{good}(\mathbf{w}^*), \mathbf{w} - \mathbf{w}^* \rangle \geq \lambda_{good} \|\mathbf{w} - \mathbf{w}^*\|^2$$

(3) (Gradient Lipschitz) $f_i(\mathbf{w})$ has L_i Lipschitz gradient i.e.,

$$\|\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}^*)\| \leq L_i \|\mathbf{w} - \mathbf{w}^*\|$$

Theorem 4. (*Distance to \mathbf{w}^**)

$$\mathbb{E}_i \left[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 | \mathbf{w}_t \right] \leq \left(1 - 2\eta_t \lambda_{good} (1 - \eta_t \sup_i L_i) \min_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \right) \|\mathbf{w}_t - \mathbf{w}^*\|^2 + R_t \quad (29)$$

where

$$\begin{aligned} R_t = & -2\eta_t \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}^*) \rangle \\ & + 2\eta_t^2 \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \|\nabla f_i(\mathbf{w}^*)\|^2 + \eta_t^2 \sum_{i \in \mathbb{O}} p_i(\mathbf{w}_t) \|\nabla f_i(\mathbf{w}_t)\|^2 + 2\eta_t \sum_{i \in \mathbb{O}} p_i(\mathbf{w}_t) (f_i(\mathbf{w}^*) - f_i(\mathbf{w}_t)) \end{aligned}$$

Proof. Observe first that for each component function i.e. ,

$$\langle \mathbf{w} - \mathbf{v}, \nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{v}) \rangle \geq \frac{1}{L_i} \|f_i(\mathbf{w}) - f_i(\mathbf{v})\|^2$$

For detailed proof, see Lemma A.1 in [Needell et al., 2014].

For each individual component function $f_i(\mathbf{w})$, we have

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 &= \|\mathbf{w}_t - \mathbf{w}^*\|^2 + \eta_t^2 \|\nabla f_i(\mathbf{w}_t)\|^2 - 2\eta_t \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) \rangle \\ &\leq \|\mathbf{w}_t - \mathbf{w}^*\|^2 + 2\eta_t^2 \|\nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*)\|^2 + 2\eta_t^2 \|\nabla f_i(\mathbf{w}^*)\|^2 - 2\eta_t \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) \rangle \\ &\leq \|\mathbf{w}_t - \mathbf{w}^*\|^2 + 2\eta_t^2 L_i \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*) \rangle + 2\eta_t^2 \|\nabla f_i(\mathbf{w}^*)\|^2 \\ &\quad - 2\eta_t \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) \rangle \\ &= \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2\eta_t (1 - \eta_t \sup_i L_i) \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*) \rangle + 2\eta_t^2 \|\nabla f_i(\mathbf{w}^*)\|^2 \\ &\quad - 2\eta_t \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}^*) \rangle \end{aligned}$$

We next take an expectation with respect to the choice of i conditional on \mathbf{w}_t

$$\begin{aligned}
 \mathbb{E}_i \left[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 | \mathbf{w}_t \right] &\leq \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2\eta_t (1 - \eta_t \sup_i L_i) \underbrace{\left\langle \mathbf{w}_t - \mathbf{w}^*, \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) (\nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*)) \right\rangle}_{Term1} \\
 &\quad - 2\eta_t \langle \mathbf{w}_t - \mathbf{w}^*, \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \nabla f_i(\mathbf{w}^*) \rangle + 2\eta_t^2 \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \|\nabla f_i(\mathbf{w}^*)\|^2 \\
 &\quad + \eta_t^2 \sum_{i \in \mathcal{O}} p_i(\mathbf{w}_t) \|\nabla f_i(\mathbf{w}_t)\|^2 + 2\eta_t \underbrace{\langle \mathbf{w}^* - \mathbf{w}_t, \sum_{i \in \mathcal{O}} p_i(\mathbf{w}_t) \nabla f_i(\mathbf{w}_t) \rangle}_{Term2}
 \end{aligned} \tag{30}$$

Now we first bound $Term1$ as follows

$$\begin{aligned}
 Term1 &\leq \min_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \sum_{i \in \mathcal{G}} \left\langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*) \right\rangle \\
 &\leq \min_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \lambda_{good} \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2
 \end{aligned}$$

For $Term2$ we apply the property of the convex function $\langle \nabla f_i(\mathbf{v}), \mathbf{w} - \mathbf{v} \rangle \leq f_i(\mathbf{w}) - f_i(\mathbf{v})$

$$Term2 \leq \sum_{i \in \mathcal{O}} p_i(\mathbf{w}_t) (f_i(\mathbf{w}^*) - f_i(\mathbf{w}_t))$$

Putting the upper bound of $Term1$ and $Term2$ back to (30) gives

$$\mathbb{E}_i \left[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 | \mathbf{w}_t \right] \leq \left(1 - 2\eta_t \lambda_{good} (1 - \eta_t \sup_i L_i) \min_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \right) \|\mathbf{w}_t - \mathbf{w}^*\|^2 + R_t \tag{31}$$

where

$$\begin{aligned}
 R_t &= -2\eta_t \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}^*) \rangle \\
 &\quad + 2\eta_t^2 \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \|\nabla f_i(\mathbf{w}^*)\|^2 + \eta_t^2 \sum_{i \in \mathcal{O}} p_i(\mathbf{w}_t) \|\nabla f_i(\mathbf{w}_t)\|^2 + 2\eta_t \sum_{i \in \mathcal{O}} p_i(\mathbf{w}_t) (f_i(\mathbf{w}^*) - f_i(\mathbf{w}_t))
 \end{aligned}$$

□

We have the following corollary that for noiseless setting, if we can have some good initialization, MKL-SGD is always better than SGD even the corrupted data is greater than half. For noisy setting, we can also perform better than SGD with one more condition: the noise is not large than the distance $\|\Delta_t\|^2$. This condition is not mild in the sense that $\|\mathbf{w}_t - \mathbf{w}^*\|^2$ is always greater than $\|\bar{\mathbf{w}}_{SGD} - \mathbf{w}^*\|^2$ for SGD algorithm and $\|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\|^2$ for MKL-SGD.

Corollary 1. *Suppose we have $|\mathcal{G}| \leq \frac{n}{2}$. At iteration t for $\eta_t \leq \frac{1}{\sup_i L_i}$, the parameter \mathbf{w}_t satisfies $\sup_{i \in \mathcal{G}} f_i(\mathbf{w}_t) \leq \inf_{j \in \mathcal{O}} f_j(\mathbf{w}_t)$. Moreover, assume the noise level at optimal \mathbf{w}^* satisfies*

$$either \quad \|\nabla f_i(\mathbf{w}^*)\| \leq \frac{\lambda_{good}(1 - \eta_t \sup_i L_i)/n}{1 + \sqrt{1 + \eta_t(1 - \eta_t \sup_i L_i)\lambda_{good}/n}} \|\mathbf{w}_t - \mathbf{w}^*\|, \text{ for } i \in \mathcal{G} \tag{32}$$

$$or \quad \sum_{i \in \mathcal{G}} \|\nabla f_i(\mathbf{w}^*)\|^2 \leq \left(\frac{\lambda_{good}(1 - \eta_t \sup_i L_i)|\mathcal{G}|/n}{\sqrt{n} + \sqrt{\sqrt{n} + \eta_t(1 - \eta_t \sup_i L_i)\lambda_{good}|\mathcal{G}|/n}} \right)^2 \|\mathbf{w}_t - \mathbf{w}^*\|^2. \tag{33}$$

Using the same setup, the vanilla SGD and MKL-SGD ($K=2$) algorithms yield respectively

$$SGD \quad \mathbb{E}_i \left[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 | \mathbf{w}_t \right] \leq \left(1 - 2\eta_t \lambda_{good} (1 - \eta_t \sup_i L_i) \frac{|\mathcal{G}|}{n} \right) \|\mathbf{w}_t - \mathbf{w}^*\|^2 + R_t^{(SGD)}$$

$$MKL-2 \quad \mathbb{E}_i \left[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 | \mathbf{w}_t \right] \leq \left(1 - 2\eta_t \lambda_{good} (1 - \eta_t \sup_i L_i) \frac{|\mathcal{G}|}{n} \right) \|\mathbf{w}_t - \mathbf{w}^*\|^2 + R_t^{(MKL_2)}$$

where

$$R_t^{(MKL_2)} \leq R_t^{(SGD)}.$$

Proof. Start from the inequality 30 in the proof of Theorem 4. We have *Term1* as follows:

$$\begin{aligned} \text{Term1} &= \frac{|\mathcal{G}|}{n} \left\langle \mathbf{w}_t - \mathbf{w}^*, \sum_{i \in \mathcal{G}} \frac{p_i(\mathbf{w}_t)}{|\mathcal{G}|/n} (\nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*)) \right\rangle \\ &= \frac{|\mathcal{G}|}{n} \langle \mathbf{w}_t - \mathbf{w}^*, \nabla F_{good}(\mathbf{w}_t) - \nabla F_{good}(\mathbf{w}^*) \rangle \\ &\quad + \frac{|\mathcal{G}|}{n} \sum_{i \in \mathcal{G}} \left(\frac{p_i(\mathbf{w}_t)}{|\mathcal{G}|/n} - \frac{1}{|\mathcal{G}|} \right) \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*) \rangle \\ &\geq \lambda_{good} \frac{|\mathcal{G}|}{n} \|\mathbf{w}_t - \mathbf{w}^*\|^2 + \sum_{i \in \mathcal{G}} \left(p_i(\mathbf{w}_t) - \frac{1}{n} \right) \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*) \rangle \end{aligned}$$

Putting the terms back to (30), we have for $\eta_t \leq 1/(\sup_i L_i)$

$$\mathbb{E}_i \left[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 | \mathbf{w}_t \right] \leq \left(1 - 2\eta_t \lambda_{good} (1 - \eta_t \sup_i L_i) \frac{|\mathcal{G}|}{n} \right) \|\mathbf{w}_t - \mathbf{w}^*\|^2 + R_t \quad (34)$$

where

$$\begin{aligned} R_t &= -2\eta_t (1 - \eta_t \sup_i L_i) \sum_{i \in \mathcal{G}} \left(p_i(\mathbf{w}_t) - \frac{1}{n} \right) \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}_t) - \nabla f_i(\mathbf{w}^*) \rangle \\ &\quad - 2\eta_t \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_i(\mathbf{w}^*) \rangle \\ &\quad + 2\eta_t^2 \sum_{i \in \mathcal{G}} p_i(\mathbf{w}_t) \|\nabla f_i(\mathbf{w}^*)\|^2 + \eta_t^2 \sum_{i \in \mathbb{O}} p_i(\mathbf{w}_t) \|\nabla f_i(\mathbf{w}_t)\|^2 + 2\eta_t \sum_{i \in \mathbb{O}} p_i(\mathbf{w}_t) (f_i(\mathbf{w}^*) - f_i(\mathbf{w}_t)) \end{aligned}$$

Now we analyse the term R_t for vanilla SGD and MKL-SGD ($K = 2$) respectively. For vanilla SGD, we have $p_i(\mathbf{w}_t) = \frac{1}{n}$ and $\sum_{i \in \mathcal{G}} \nabla f_i(\mathbf{w}^*) = 0$, which results in

$$R_t^{(SGD)} = \frac{2\eta_t^2}{n} \sum_{i \in \mathcal{G}} \|\nabla f_i(\mathbf{w}^*)\|^2 + \frac{\eta_t^2}{n} \sum_{i \in \mathbb{O}} \|\nabla f_i(\mathbf{w}_t)\|^2 + \frac{2\eta_t}{n} \sum_{i \in \mathbb{O}} (f_i(\mathbf{w}^*) - f_i(\mathbf{w}_t))$$

Note that MKL-SGD for $K = 2$ have

$$p_{m_i(\mathbf{w})}(\mathbf{w}) = \frac{2(n-i)}{n(n-1)} \quad (35)$$

where $m_1(\mathbf{w}), m_2(\mathbf{w}), m_3(\mathbf{w}), \dots, m_n(\mathbf{w})$ are the indices of data samples for some \mathbf{w} :

$$f_{m_1(\mathbf{w})}(\mathbf{w}) \leq f_{m_2(\mathbf{w})}(\mathbf{w}) \leq \dots \leq f_{m_n(\mathbf{w})}(\mathbf{w})$$

Suppose the iteration \mathbf{w}_t satisfies that $f_i(\mathbf{w}_t) < f_j(\mathbf{w}_t)$ for $i \in \mathcal{G}, j \in \mathbb{O}$. For $|\mathcal{G}| \leq \frac{n}{2}$, we have for

$$\begin{aligned}
 R_t^{(MKL_2)} &= -2\eta_t(1 - \eta_t \sup_i L_i) \sum_{i=1}^{|\mathcal{G}|} \frac{(n-2i+1)}{n(n-1)} \langle \mathbf{w}_t - \mathbf{w}^*, \nabla f_{m_i}(\mathbf{w}_t) - \nabla f_{m_i}(\mathbf{w}^*) \rangle \\
 &\quad + 2\eta_t \sum_{i=1}^{|\mathcal{G}|} \frac{2(n-i)}{n(n-1)} \left(\langle \mathbf{w}^* - \mathbf{w}_t, \nabla f_i(\mathbf{w}^*) \rangle + \eta_t \|\nabla f_i(\mathbf{w}^*)\|^2 \right) \\
 &\quad + \eta_t^2 \sum_{i=|\mathcal{G}|+1}^n \frac{2(n-i)}{n(n-1)} \|\nabla f_i(\mathbf{w}_t)\|^2 + 2\eta_t \sum_{i=|\mathcal{G}|+1}^n \frac{2(n-i)}{n(n-1)} (f_i(\mathbf{w}^*) - f_i(\mathbf{w}_t)) \\
 &\leq -2\eta_t(1 - \eta_t \sup_i L_i) \frac{|\mathcal{G}| \lambda_{good}}{n(n-1)} \|\mathbf{w}_t - \mathbf{w}^*\|^2 \\
 &\quad + \frac{4\eta_t}{n} \sum_{i=1}^{|\mathcal{G}|} \left(\|\mathbf{w}^* - \mathbf{w}_t\| \|\nabla f_i(\mathbf{w}^*)\| + \eta_t \|\nabla f_i(\mathbf{w}^*)\|^2 \right) \\
 &\quad + \sum_{i=|\mathcal{G}|+1}^n \frac{\eta_t^2}{n} \|\nabla f_i(\mathbf{w}_t)\|^2 + \sum_{i=|\mathcal{G}|+1}^n \frac{2\eta_t}{n} (f_i(\mathbf{w}^*) - f_i(\mathbf{w}_t))
 \end{aligned}$$

We will have $R_t^{(MKL_2)} \leq R_t^{(SGD)}$ if the following inequality holds

$$(1 - \eta_t \sup_i L_i) \frac{|\mathcal{G}| \lambda_{good}}{(n-1)} \|\mathbf{w}_t - \mathbf{w}^*\|^2 \geq \sum_{i=1}^{|\mathcal{G}|} \left(2\|\mathbf{w}^* - \mathbf{w}_t\| \|\nabla f_i(\mathbf{w}^*)\| + \eta_t \|\nabla f_i(\mathbf{w}^*)\|^2 \right). \quad (36)$$

Indeed, for the noise level $\|\nabla f_i(\mathbf{w}^*)\|^2$ satisfying (32) we have for $i \in \mathcal{G}$,

$$(1 - \eta_t \sup_i L_i) \frac{\lambda_{good}}{(n-1)} \|\mathbf{w}_t - \mathbf{w}^*\|^2 \geq 2\|\mathbf{w}^* - \mathbf{w}_t\| \|\nabla f_i(\mathbf{w}^*)\| + \eta_t \|\nabla f_i(\mathbf{w}^*)\|^2.$$

Summing up the terms in $i \in \mathcal{G}$, we get (36). For the noise level $\|\nabla f_i(\mathbf{w}^*)\|^2$ satisfying (33) we have

$$\begin{aligned}
 (1 - \eta_t \sup_i L_i) \frac{\lambda_{good} |\mathcal{G}|}{(n-1)} \|\mathbf{w}_t - \mathbf{w}^*\|^2 &\geq \left(2\|\mathbf{w}^* - \mathbf{w}_t\| \sqrt{n \sum_{i \in \mathcal{G}} \|\nabla f_i(\mathbf{w}^*)\|^2} + \eta_t \sum_{i \in \mathcal{G}} \|\nabla f_i(\mathbf{w}^*)\|^2 \right) \\
 &\geq 2\|\mathbf{w}^* - \mathbf{w}_t\| \sum_{i \in \mathcal{G}} \|\nabla f_i(\mathbf{w}^*)\| + \eta_t \sum_{i \in \mathcal{G}} \|\nabla f_i(\mathbf{w}^*)\|^2.
 \end{aligned}$$

which results in (36). \square

7.4 More experimental results

7.4.1 Linear Regression

Here, we show that there exists a trade-off for MKL-SGD between the rate of convergence and robustness the algorithm provides against outliers depending on the value of the parameter k . Larger the k , more robust is the algorithm, but slower is the rate of convergence. The algorithm outperforms median loss SGD and SGD. We also experimented with other order statistics and observed that for most general settings MKL-SGD was the best to pick. Note that the outliers are chosen from $\mathcal{N}(0, 1)$ distribution independently of the data sample.

Choosing the Sample with Lowest Loss makes SGD Robust

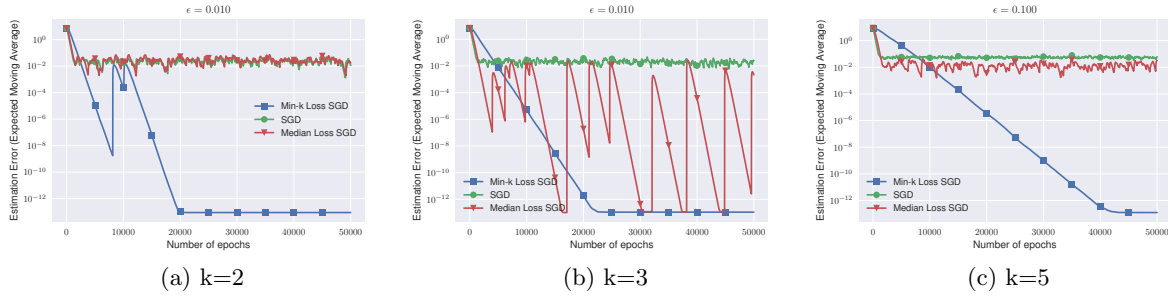


Figure 6: Comparing the performance of MKL-SGD , SGD and Median loss SGD in the noiseless setting, $d = 50$.

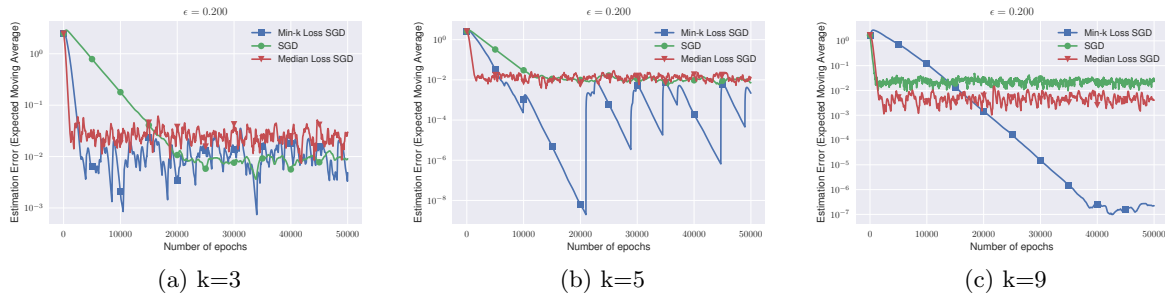


Figure 7: Comparing the performance of MKL-SGD , SGD and Median loss SGD in the noisy setting, $d = 10$, Noise variance=0.0001

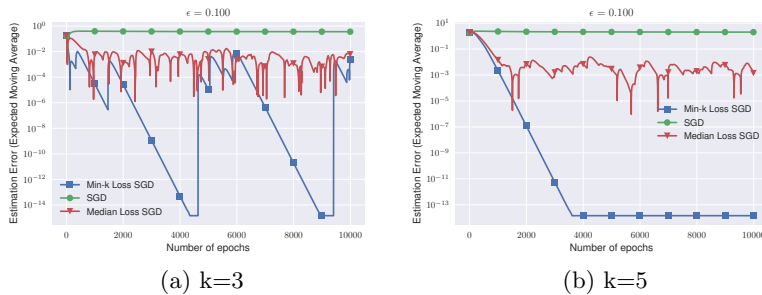


Figure 8: Comparing the performance of MKL-SGD , SGD and Median loss SGD in the noiseless setting, $d = 25$, Noise variance=0.01

7.4.2 Neural Network Experiments

Here, we show that in presence of outliers instead of tuning other hyperparameters like learning rate, tuning over k might lead to significant gains in performances for deep neural networks. To illustrate this we play around with two commonly used noise models: random noise and directed noise. In the random noise model, the outlier label is randomly assigned while for the directed noise model for some class ‘a’, the outlier is assigned the same label ‘b’, similarly all the outliers for class ‘b’ are assigned label ‘c’ and so on.

Dataset	MNIST with 2-layer CNN (Directed Noise)						
Optimizer	SGD	MKL-SGD					Oracle
$\epsilon \backslash \alpha$	1.0	0.9	0.8	0.7	0.6	0.5	1.0
0.1	96.76	97.23	95.89	97.47	96.34	94.54	98.52
0.2	92.54	95.81	95.58	97.46	97.03	95.76	98.33
0.3	85.77	91.56	93.59	95.30	96.54	95.96	98.16
0.4	71.95	78.68	82.25	85.93	91.29	94.20	97.98

Table 2: In this experiments, we train a standard 2 layer CNN on subsampled MNIST (5000 training samples with labels corrupted using random label noise). We train over 80 epochs using an initial learning rate of 0.05 with the decaying schedule of factor 5 after every 30 epochs. The reported accuracy is based on the true validation set. The results of the MNIST dataset are reported as the mean of 5 runs. For the MKL-SGD algorithm, we introduce a more practical variant that evaluates k sample losses and picks a batch of size αk where $k = 10$.

Dataset	MNIST with 2-layer CNN (Random Noise)						
Optimizer	SGD	MKL-SGD					Oracle
$\epsilon \backslash \alpha$	1.0	0.9	0.8	0.7	0.6	0.5	1.0
0.1	96.91	97.9	98.06	97.59	96.49	94.43	98.44
0.2	93.94	95.5	96.16	97.02	97.04	96.25	98.18
0.3	87.14	90.71	91.60	92.97	94.54	95.36	97.8
0.4	71.83	74.31	76.6	78.30	77.58	80.86	97.16

Table 3: In this experiments, we train a standard 2 layer CNN on subsampled MNIST (5000 training samples with labels corrupted using random label noise). We train over 80 epochs using an initial learning rate of 0.05 with the decaying schedule of factor 5 after every 30 epochs. The reported accuracy is based on the true validation set. The results of the MNIST dataset are reported as the mean of 5 runs. For the MKL-SGD algorithm, we introduce a more practical variant that evaluates k sample losses and picks a batch of size αk where $k = 10$.

Dataset	CIFAR-10 with Resnet-18 (Directed Noise)						
Optimizer	SGD	MKL-SGD					Oracle
$\epsilon \backslash \alpha$	1.0	0.9	0.8	0.7	0.6	0.5	1.0
0.1	79.1	77.52	79.57	81.00	81.94	80.53	84.56
0.2	72.29	69.58	70.17	72.76	77.77	78.93	84.40
0.3	63.96	61.43	60.46	61.58	66.49	69.57	84.66
0.4	52.4	51.53	51.04	51.07	53.57	51.2	84.42

Table 4: In this experiments, we train Resnet 18 on CIFAR-10 (50000 training samples with labels corrupted using directed label noise). We train over 200 epochs using an initial learning rate of 0.05 with the decaying schedule of factor 5 after every 90 epochs. The reported accuracy is based on the true validation set. The results of the CIFAR-10 dataset are reported as the mean of 3 runs. For the MKL-SGD algorithm, we introduce a more practical variant that evaluates k sample losses and picks a batch of size αk where $k = 16$.

7.5 Conclusions and Future Work

7.5.1 Conclusions

In this paper, we propose MKL-SGD that is computationally inexpensive, has linear convergence (upto a certain neighborhood) and is robust against outliers. We analyze MKL-SGD algorithm under noiseless and noisy settings with and without outliers. MKL-SGD outperforms SGD in terms of generalization for both

linear regression and neural network experiments. MKL-SGD opens up a plethora of challenging questions with respect to understanding convex optimization in a non-convex landscape which will be discussed in the Appendix.

7.5.2 Future Work

To ensure consistency, i.e. $\|\bar{\mathbf{w}}_{MKL} - \mathbf{w}^*\| \rightarrow 0$, we require that $k \geq n\epsilon + 1$. In all other cases, there will be a non-zero contribution from the outliers which keeps the MKL-SGD solution from exactly converging to \mathbf{w}^* . In this paper, we consider unknown ϵ and thus k should be a hyperparameter. For neural network experiments in the Appendix, we show that tuning k as a hyperparameter can lead to significant improvements in performance in presence of outliers.

The obvious question is if it is possible to provide worst case guarantees for a larger subset of problems using smarter initialization techniques. It will be interesting to analyze the tradeoff between better generalization guarantees offered by large k and rates of convergence. The worst case analysis in the noisy setting for standard convex optimization losses remains an open problem. As we show in the previous set of experiments, in presence of noise, tuning the hyperparameter k can provide significant boosts to the performance.