# TRAINING DEEP NEURAL NETWORKS VIA DIRECT LOSS MINIMIZATION: SUPPLEMENTARY MATERIALS

## 1 PROOF OF THE GENERAL LOSS GRADIENT THEOREM

In order to lay the foundation for the proof of the general loss gradient theorem, we first show the following lemma. In short, it provides the bases for exchanging integral bounds when $\epsilon$ approaches $0$ from above. All integrals used in this section should be viewed as Lebesgue integrals.

**Lemma 1.**

$$\lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_{b\epsilon+o(\epsilon)}^{\infty} f(x,y,\epsilon)dxdy = \lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^{\infty} f(x,y,\epsilon)dxdy$$

*Proof.* We have

$$\lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_{b\epsilon+o(\epsilon)}^{\infty} f(x,y,\epsilon)dxdy$$

$$= \lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^{\infty} f(x,y,\epsilon)dxdy - \lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^{b\epsilon+o(\epsilon)} f(x,y,\epsilon)dxdy.$$

Suppose $f$ is continuous w.r.t $x, y, \epsilon$, then as $\epsilon \to 0^+$, it can be bounded by some constant $M$. As a result, we have

$$\frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^{b\epsilon+o(\epsilon)} |f(x,y,\epsilon)|dxdy \leq M(ab\epsilon + ao(\epsilon) + bo(\epsilon) + o(\epsilon)o(\epsilon)/\epsilon), \quad (1)$$

which means

$$\lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^{b\epsilon+o(\epsilon)} f(x,y,\epsilon)dxdy = 0. \quad (2)$$

$\square$

Repeated application of Lemma 1 as demonstrated in the following is directly helpful for the proof of the general loss gradient theorem, which is why we state it explicitly:

**Lemma 2.** *Let $a > 0$, then we assert*

$$\lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_{b_1\epsilon+o(\epsilon)}^{\infty} \cdots \int_{b_n\epsilon+o(\epsilon)}^{\infty} f(x, y_1, \cdots, y_n)dxdy_1 \cdots dy_n$$

$$= \lim_{\epsilon \to 0^+} \frac{1}{\epsilon} a \int_0^{\epsilon} \int_0^{\infty} \cdots \int_0^{\infty} f(x, y_1, \cdots, y_n)dxdy_1 \cdots dy_n$$

*Proof.* Let $f(x, y_1, \epsilon) = \int_{b_2\epsilon+o(\epsilon)}^{\infty} \cdots \int_{b_n\epsilon+o(\epsilon)}^{\infty} f(x, y_1, \cdots, y_n)dy_2 \cdots dy_n$. Due to Lemma 1 we obtain the following:

$$\lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_{b_1\epsilon+o(\epsilon)}^{\infty} f(x, y_1, \epsilon)dxdy_1 = \lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^{\infty} f(x, y_1, \epsilon)dxdy_1$$

$$= \lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^{\infty} \int_{b_2\epsilon+o(\epsilon)}^{\infty} \cdots \int_{b_n\epsilon+o(\epsilon)}^{\infty} f(x, y_1, \cdots, y_n)dxdy_1 \cdots dy_n.$$

Now we denote $f(x, y_2, \epsilon) = \int_0^\infty \int_{b_3\epsilon+o(\epsilon)}^\infty \cdots \int_{b_n\epsilon+o(\epsilon)}^\infty f(x, y_1, \cdots, y_n) dy_1 dy_3 \cdots dy_n$ and follow a similar procedure:

$$
\lim_{\epsilon\to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^\infty \int_{b_2\epsilon+o(\epsilon)}^\infty \cdots \int_{b_n\epsilon+o(\epsilon)}^\infty f(x, y_1, \cdots, y_n) dx dy_1 \cdots dy_n
$$

$$
= \lim_{\epsilon\to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_{b_2\epsilon+o(\epsilon)}^\infty f(x, y_2, \epsilon) dx dy_2
$$

$$
= \lim_{\epsilon\to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^\infty f(x, y_2, \epsilon) dx dy_2
$$

$$
= \cdots \text{(following this procedure recursively)}
$$

$$
= \lim_{\epsilon\to 0^+} \frac{1}{\epsilon} \int_0^{a\epsilon+o(\epsilon)} \int_0^\infty \cdots \int_0^\infty f(x, y_1, \cdots, y_n) dx dy_1 \cdots dy_n
$$

$$
= \lim_{\epsilon\to 0^+} \frac{a\epsilon + o(\epsilon)}{\epsilon} \lim_{\epsilon\to 0^+} \frac{1}{a\epsilon + o(\epsilon)} \int_0^{a\epsilon+o(\epsilon)} \int_0^\infty \cdots \int_0^\infty f(x, y_1, \cdots, y_n) dx dy_1 \cdots dy_n
$$

$$
= a \lim_{\epsilon\to 0^+} \frac{1}{a\epsilon + o(\epsilon)} \int_0^{a\epsilon+o(\epsilon)} \int_0^\infty \cdots \int_0^\infty f(x, y_1, \cdots, y_n) dx dy_1 \cdots dy_n
$$

$$
= a \lim_{a\epsilon+o(\epsilon)\to 0^+} \frac{1}{a\epsilon + o(\epsilon)} \int_0^{a\epsilon+o(\epsilon)} \int_0^\infty \cdots \int_0^\infty f(x, y_1, \cdots, y_n) dx dy_1 \cdots dy_n
$$

$$
= \lim_{\epsilon\to 0^+} \frac{1}{\epsilon} a \int_0^\epsilon \int_0^\infty \cdots \int_0^\infty f(x, y_1, \cdots, y_n) dx dy_1 \cdots dy_n.
$$

This completes the proof. □

For readability we repeat the main theorem:

**Theorem 1** (General Loss Gradient Theorem). *When given a finite set $\mathcal{Y}$, a scoring function $F(x, y, w)$, a data distribution, as well as a task-loss $L(y, \hat{y})$, then, under some mild regularity conditions (see the proof for details), the direct loss gradient has the following form:*

$$
\nabla_w \mathbb{E}\left[L(y, y_w)\right] = \lim_{\epsilon\to 0} \frac{\pm 1}{\epsilon} \mathbb{E}\left[\nabla_w F(x, y_{\text{direct}}, w) - \nabla_w F(x, y_w, w)\right],
$$

*with*

$$
y_w = \arg\max_{\hat{y}\in\mathcal{Y}} F(x, \hat{y}, w),
$$

$$
y_{\text{direct}} = \arg\max_{\hat{y}\in\mathcal{Y}} F(x, \hat{y}, w) \pm \epsilon L(y, \hat{y}).
$$

In the following we prove the positive case and note that the negative case is easily proved using a similar procedure.

*Proof.* Without loss of generality, in this proof we assume $\mathcal{Y} = \{1, 2, \ldots, |\mathcal{Y}|\}$ and no tie in maximization.

By definition of the directional derivative we have

$$
\Delta w^\intercal \nabla_w \mathbb{E}\left[L(y, y_w(x)\right] = \lim_{\epsilon\to 0} \frac{\mathbb{E}\left[L(y, y_{w+\epsilon\Delta w}(x))\right] - \mathbb{E}\left[L(y, y_w(x))\right]}{\epsilon}.
$$

Hence we need to prove the following equivalence:

$$
\lim_{\epsilon\to 0} \frac{\mathbb{E}\left[L(y, y_{w+\epsilon\Delta w}(x)) - L(y, y_w(x))\right]}{\epsilon} \tag{3}
$$

$$
= \lim_{\epsilon\to 0} \frac{\Delta w^\intercal \mathbb{E}\left[\nabla_w F(x, y_{\text{direct}}, w) - \nabla_w F(x, y_w(x), w)\right]}{\epsilon}. \tag{4}
$$

Denote $\Delta F_w^{i,j}(x) = F(x,i,w) - F(x,j,w)$, $\Delta L(y)^{i,j} = L(y,i) - L(y,j)$. Note that $\Delta F_w^{i,i}(x) = 0$ and $\Delta L(y)^{i,i} = 0$. Therefore we just have to consider terms where the classification result changes when moving from $w$ to $w + \epsilon \Delta w$. To this end we decompose the expectation in Eq. (3) into pairwise terms to yield

$$\mathbb{E}\left[L(y, y_{w+\epsilon\Delta w}(x) - L(y, y_w(x)))\right] = \sum_{i \neq j} \mathbb{E}\left[\Delta L(y)^{i,j} \mathbf{1}_{\{x \in \mathcal{A}^{i,j}\}}\right],$$

where the indicator set for a change from class label $i$ to category $j$ when moving from $w$ to $w + \Delta w$ is given by

$$\begin{aligned}
\mathcal{A}^{i,j} &= \{x : y_{w+\Delta w}(x) = i, y_w(x) = j\} \\
&= \{x : \Delta F_{w+\epsilon\Delta w}^{i,k} > 0, \Delta F_w^{j,k} > 0, \quad \forall k \in \mathcal{Y}\} \\
&= \{x : \Delta F_w^{i,k} + \epsilon\Delta w^\mathsf{T} \nabla \Delta F_w^{i,k} + o(\epsilon) > 0, \Delta F_w^{j,k} > 0, \quad \forall k \in \mathcal{Y}\} \\
&= \{x : 0 < \Delta F_w^{j,i} < -\epsilon\Delta w^\mathsf{T} \nabla \Delta F_w^{j,i} + o(\epsilon) \\
&\quad \Delta F_w^{j,k} > 0, \quad k \neq i \\
&\quad \Delta F_w^{i,k} > -\epsilon\Delta w^\mathsf{T} \nabla \Delta F_w^{i,k} + o(\epsilon), \quad k \neq j\}.
\end{aligned}$$

As in (McAllester et al., 2010), we assume that any joint measure $\rho$ on $\Delta F_w^{j,1} \cdots \Delta F_w^{i,n}, \Delta L(y)^{i,j}, \Delta w^\mathsf{T} \nabla \Delta F_w^{i,j}$ can be expressed as a measure $\mu$ on $\Delta L(y)^{i,j}, \Delta w^\mathsf{T} \nabla \Delta F_w^{i,j}$ and a bounded continuous conditional density function $f$. Integrating the loss difference over the set, we obtain

$$\begin{aligned}
&\mathbb{E}\left[\Delta L(y)^{i,j} \mathbf{1}_{\{x \in \mathcal{A}^{i,j}\}}\right] \\
&= \mathbb{E}_\mu\Bigg[\Delta L(y)^{i,j} \int_0^{\epsilon\Delta w^\mathsf{T} \nabla \Delta F_w^{i,j} + o(\epsilon)} d\Delta F_w^{j,i} \int_0^\infty d\Delta F_w^{j,1} \cdots \int_0^\infty d\Delta F_w^{j,n} \\
&\int_{-\epsilon\Delta w^\mathsf{T} \nabla \Delta F_w^{i,1} + o(\epsilon)}^\infty d\Delta F_w^{i,1} \cdots \\
&\int_{-\epsilon\Delta w^\mathsf{T} \nabla \Delta F_w^{i,n} + o(\epsilon)}^\infty f(\Delta F_w^{j,1} \cdots \Delta F_w^{i,n} \mid \Delta L(y)^{i,j}, \Delta w^\mathsf{T} \nabla \Delta F_w^{i,j}) d\Delta F_w^{i,n}\Bigg].
\end{aligned}$$

Based on Lemma 2, we conclude that Eq. (3) is equivalent to

$$\begin{aligned}
&\lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \mathbb{E}\left[\Delta L(y)^{i,j} \mathbf{1}_{\{x \in \mathcal{A}^{i,j}\}}\right] \\
&= \mathbb{E}_\mu\Bigg[\Delta L(y)^{i,j} (\Delta w^\mathsf{T} \nabla \Delta F_w^{i,j})^+ \lim_{\epsilon \to 0^+} \frac{1}{\epsilon} \int_0^\epsilon d\Delta F_w^{j,i} \int_0^\infty d\Delta F_w^{j,1} \cdots \int_0^\infty d\Delta F_w^{j,n} \\
&\int_0^\infty d\Delta F_w^{i,1} \cdots \int_0^\infty f(\Delta F_w^{j,1} \cdots \Delta F_w^{i,n} | \Delta L(y)^{i,j}, \Delta w^\mathsf{T} \nabla \Delta F_w^{i,j}) d\Delta F_w^{i,n}\Bigg].
\end{aligned} \quad (5)$$

Following a similar procedure, we decompose the expectation in Eq. (4) to

$$\mathbb{E}\left[\Delta w^\mathsf{T} \nabla F(x, y_{\text{direct}}, w) - \Delta w^\mathsf{T} \nabla F(x, y_w(x), w)\right] = \sum_{i \neq j} \mathbb{E}\left[\Delta w^\mathsf{T} \Delta \nabla F(x)^{i,j} \mathbf{1}_{\{(x,y) \in \mathcal{B}^{i,j}\}}\right],$$

where the indicator set for a change from label $i$ to a configuration $j$ when changing from $w$ to loss augmented inference is given by

$$\begin{aligned}
\mathcal{B}^{i,j} &= \{(x,y) : y_{\text{direct}} = i, y_w(x) = j\} \\
&= \{(x,y) : \Delta F_w^{i,k} + \epsilon\Delta L^{i,k} > 0, \Delta F_w^{j,k} > 0, \quad k \in \mathcal{Y}\} \\
&= \{(x,y) : 0 < \Delta F_w^{j,i} < -\epsilon\Delta L(y)^{j,i} \\
&\quad \Delta F_w^{j,k} > 0, \quad k \neq i \\
&\quad \Delta F_w^{i,k} > -\epsilon\Delta L(y)^{i,k}, \quad k \neq j\}.
\end{aligned}$$

Integrating the directional derivative over the set of configuration changes, we obtain

$$
\mathbb{E}\left[\Delta w^{\mathsf{T}}\Delta\nabla F(x)^{i,j}\mathbf{1}_{\{(x,y)\in\mathcal{B}^{i,j}\}}\right]
$$

$$
=\mathbb{E}_{\mu}\left[\Delta w^{\mathsf{T}}\Delta\nabla F(x)^{i,j}\int_{0}^{\epsilon\Delta L(y)^{i,j}}d\Delta F_{w}^{j,i}\int_{0}^{\infty}d\Delta F_{w}^{j,1}\cdots\int_{0}^{\infty}d\Delta F_{w}^{j,n}\right.
$$

$$
\int_{-\epsilon\Delta L(y)^{i,1}}^{\infty}d\Delta F_{w}^{i,1}\cdots
$$

$$
\left.\int_{-\epsilon\Delta L(y)^{i,n}}^{\infty}f(\Delta F_{w}^{j,1}\cdots\Delta F_{w}^{i,n}|\Delta w^{\mathsf{T}}\Delta\nabla F_{w}(x)^{i,j},\Delta L(y)^{j,i})d\Delta F_{w}^{i,n}\right].
$$

Assuming bounded data distribution and a bounded and continuous integrand we can exchange the limit operation and expectation to get

$$
\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\mathbb{E}\left[\Delta w^{\mathsf{T}}\Delta\nabla F(x)^{i,j}\mathbf{1}_{\{(x,y)\in\mathcal{B}^{i,j}\}}\right]
$$

$$
=\mathbb{E}_{\mu}\left[(\Delta L(y)^{i,j})^{+}\Delta w^{\mathsf{T}}\Delta\nabla F_{w}^{i,j}\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\int_{0}^{\epsilon}d\Delta F_{w}^{j,i}\int_{0}^{\infty}d\Delta F_{w}^{j,1}\cdots\int_{0}^{\infty}d\Delta F_{w}^{j,n}\right. \tag{6}
$$

$$
\left.\int_{0}^{\infty}d\Delta F_{w}^{i,1}\cdots\int_{0}^{\infty}f(\Delta F_{w}^{j,1}\cdots\Delta F_{w}^{i,n}|\Delta w^{\mathsf{T}}\Delta\nabla F_{w}(x)^{i,j},\Delta L(y)^{j,i})d\Delta F_{w}^{i,n}\right].
$$

Next we group expectations for a change from label $i$ to configuration $j$ and the reverse. To this end we first consider the resulting Eq. (5) obtained from rephrasing Eq. (3). Combining both label change directions yields

$$
\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\mathbb{E}\left[\Delta L(y)^{i,j}\mathbf{1}_{\{x\in\mathcal{A}^{i,j}\}}\right]+\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\mathbb{E}\left[\Delta L(y)^{j,i}\mathbf{1}_{\{x\in\mathcal{A}^{j,i}\}}\right]
$$

$$
=\mathbb{E}_{\mu}\left[\Delta L(y)^{i,j}(\Delta w^{\mathsf{T}}\nabla\Delta F_{w}^{i,j})^{+}\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\int_{0}^{\epsilon}d\Delta F_{w}^{j,i}\int_{0}^{\infty}d\Delta F_{w}^{j,1}\cdots\int_{0}^{\infty}d\Delta F_{w}^{j,n}\right.
$$

$$
\left.\int_{0}^{\infty}d\Delta F_{w}^{i,1}\cdots\int_{0}^{\infty}f(\Delta F_{w}^{j,1}\cdots\Delta F_{w}^{i,n}|\Delta L(y)^{i,j},\Delta w^{\mathsf{T}}\nabla\Delta F_{w}^{i,j})d\Delta F_{w}^{i,n}\right]
$$

$$
+\mathbb{E}_{\mu}\left[\Delta L(y)^{j,i}(\Delta w^{\mathsf{T}}\nabla\Delta F_{w}^{j,i})^{+}\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\int_{0}^{\epsilon}d\Delta F_{w}^{i,j}\int_{0}^{\infty}d\Delta F_{w}^{i,1}\cdots\int_{0}^{\infty}d\Delta F_{w}^{i,n}\right. \tag{7}
$$

$$
\left.\int_{0}^{\infty}d\Delta F_{w}^{j,1}\cdots\int_{0}^{\infty}f(\Delta F_{w}^{j,1}\cdots\Delta F_{w}^{i,n}|\Delta L(y)^{i,j},\Delta w^{\mathsf{T}}\nabla\Delta F_{w}^{i,j})d\Delta F_{w}^{i,n}\right]
$$

$$
=\mathbb{E}_{\mu}\left[\Delta L(y)^{i,j}\Delta w^{\mathsf{T}}\nabla\Delta F_{w}^{i,j}\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\int_{0}^{\epsilon}d\Delta F_{w}^{j,i}\int_{0}^{\infty}d\Delta F_{w}^{j,1}\cdots\int_{0}^{\infty}d\Delta F_{w}^{j,n}\right.
$$

$$
\left.\int_{0}^{\infty}d\Delta F_{w}^{i,1}\cdots\int_{0}^{\infty}f(\Delta F_{w}^{j,1}\cdots\Delta F_{w}^{i,n}|\Delta L(y)^{i,j},\Delta w^{\mathsf{T}}\nabla\Delta F_{w}^{i,j})d\Delta F_{w}^{i,n}\right].
$$

Similarly, we group expectations for both label change directions for the resulting Eq. (6) obtained from rephrasing Eq. (4), which yields

$$
\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\mathbb{E}\left[\Delta w^{\mathsf{T}}\Delta\nabla F(x)^{i,j}\mathbf{1}_{\{(x,y)\in\mathcal{B}^{i,j}\}}\right]+\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\mathbb{E}\left[\Delta w^{\mathsf{T}}\Delta\nabla F(x)^{j,i}\mathbf{1}_{\{(x,y)\in\mathcal{B}^{j,i}\}}\right]
$$

$$
=\mathbb{E}_{\mu}\left[\Delta L(y)^{i,j}\Delta w^{\mathsf{T}}\Delta\nabla F_{w}^{i,j}\lim_{\epsilon\to0^{+}}\frac{1}{\epsilon}\int_{0}^{\epsilon}d\Delta F_{w}^{j,i}\int_{0}^{\infty}d\Delta F_{w}^{j,1}\cdots\int_{0}^{\infty}d\Delta F_{w}^{j,n}\right. \tag{8}
$$

$$
\left.\int_{0}^{\infty}d\Delta F_{w}^{i,1}\cdots\int_{0}^{\infty}f(\Delta F_{w}^{j,1}\cdots\Delta F_{w}^{i,n}|\Delta w^{\mathsf{T}}\Delta\nabla F_{w}(x)^{i,j},\Delta L(y)^{j,i})d\Delta F_{w}^{i,n}\right].
$$

We therefore have equivalence between Eq. (7) and Eq. (8) for a change from configuration $i$ to $j$ and the reverse. Since this holds for all pairwise configurations, Eq. (3) is identical to Eq. (4), which proves the theorem. $\qquad\square$

The conditions for the above results to hold are similar to the conditions for the proof for the binary linear case (McAllester et al., 2010). The conditions can be inferred from the proof above. We

require that the joint measure $\rho$ can be expressed as a measure $\mu$ and a corresponding bounded continuous conditional density function $f$. For exchangeability of limits and expectations, it is sufficient to require the integrand to be continuous and bounded as well as the range of integral to be bounded, *i.e.*, the range of data is bounded. Further we require the scoring function $F$ to have continuous derivatives w.r.t. $w$.

## 2   PROOF FOR LEMMA 1

Next we provide the proof for Lemma 1 which we repeat for completeness. We note again that the cost function value obtained when restricting loss-augmented inference to the subsets can be computed as:

$$h(i,j) = \max_{\hat{y}} \frac{1}{|\mathcal{P}||\mathcal{N}|} \sum_{m \in \mathcal{P}_i} \sum_{n \in \mathcal{N}_j} \hat{y}_{m,n}(\phi(x_m, w) - \phi(x_n, w)) \pm \epsilon L_{AP}^{i,j}(\text{rank}(y), \text{rank}(\hat{y})), \quad (9)$$

where $L_{AP}^{i,j}$ refers to the AP loss restricted to subsets of $i$ positive and $j$ negative elements.

**Lemma 1.** *Suppose that* $\text{rank}(\hat{y}^*)$ *is the optimal ranking for Eq. (9) when restricted to $i$ positive and $j$ negative samples. Any of its sub-sequences starting at position 1 is then also an optimal ranking for the corresponding restricted sub-problem.*

*Proof.* We consider the prefix $r_1, r_2, \cdots, r_k$, where $k < |\mathcal{P}| + |\mathcal{N}|$ and $(r_1, r_2, \cdots, r_{|\mathcal{P}|+|\mathcal{N}|}) :=$ $\text{rank}(\hat{y}^*)$. Suppose there are $i$ relevant objects and $j$ irrelevant objects in the prefix, and $k = i + j$. What we need to prove is that $r_1, r_2, \cdots, r_{i+j}$ is already an optimal ranking.

Let

$$h(i,j) = \max_{\hat{y}} \underbrace{\frac{1}{|\mathcal{P}||\mathcal{N}|} \sum_{m \in \mathcal{P}_i} \sum_{n \in \mathcal{N}_j} \hat{y}_{m,n}(\phi(x_m, w) - \phi(x_n, w)) \pm \epsilon L_{AP}^{i,j}(\text{rank}(y), \text{rank}(\hat{y}))}_{:=s(i,j)}.$$

We decompose the optimal value obtained when considering all samples, *i.e.*, $h(|\mathcal{P}|, |\mathcal{N}|)$, into three parts:

$$h(|\mathcal{P}|, |\mathcal{N}|)$$
$$= \underbrace{\frac{1}{|\mathcal{P}||\mathcal{N}|} \sum_{m \in \mathcal{P}_i} \sum_{n \in \mathcal{N}_j} \hat{y}_{m,n}^*(\phi(x_m, w) - \phi(x_n, w)) \pm \epsilon L_{AP}^{i,j}(\text{rank}(y), \text{rank}(\hat{y}^*))}_{\text{Prefix terms} = s(i,j)}$$
$$+ \underbrace{\frac{1}{|\mathcal{P}||\mathcal{N}|} \sum_{m \in \mathcal{P} \setminus \mathcal{P}_i} \sum_{n \in \mathcal{N} \setminus \mathcal{N}_j} \hat{y}_{m,n}^*(\phi(x_m, w) - \phi(x_n, w)) \pm \epsilon L_{AP}^{\neg i, \neg j}(\text{rank}(y), \text{rank}(\hat{y}^*))}_{\text{Suffix terms}}$$
$$+ \underbrace{\frac{1}{|\mathcal{P}||\mathcal{N}|} \left[ \sum_{m \in \mathcal{P}_i} \sum_{n \in \mathcal{N} \setminus \mathcal{N}_j} \hat{y}_{m,n}^*(\phi(x_m, w) - \phi(x_n, w)) + \sum_{m \in \mathcal{P} \setminus \mathcal{P}_i} \sum_{n \in \mathcal{N}_j} \hat{y}_{m,n}^*(\phi(x_m, w) - \phi(x_n, w)) \right]}_{\text{Cross terms}}.$$

Here, $L_{AP}^{\neg i, \neg j}(\text{rank}(y), \text{rank}(\hat{y}^*))$ refers to the loss obtained by considering all the samples not within the prefix.

Intuitively, when changing the interleaving pattern of the prefix, the suffix terms and cross terms remain the same. This is true since the suffix terms are independent of the ranking of the prefix terms. In addition the cross terms only depend on the number and scores of positive and negative elements in the prefix but not their specific ranking.

More formally, suppose that $f(i,j) \neq h(i,j)$, then we can substitute $h(i,j)$ into the prefix term and get a larger value than $h(|\mathcal{P}|, |\mathcal{N}|)$, contradicting the fact that $h(|\mathcal{P}|, |\mathcal{N}|)$ is already the largest, which concludes the proof. $\qquad \square$
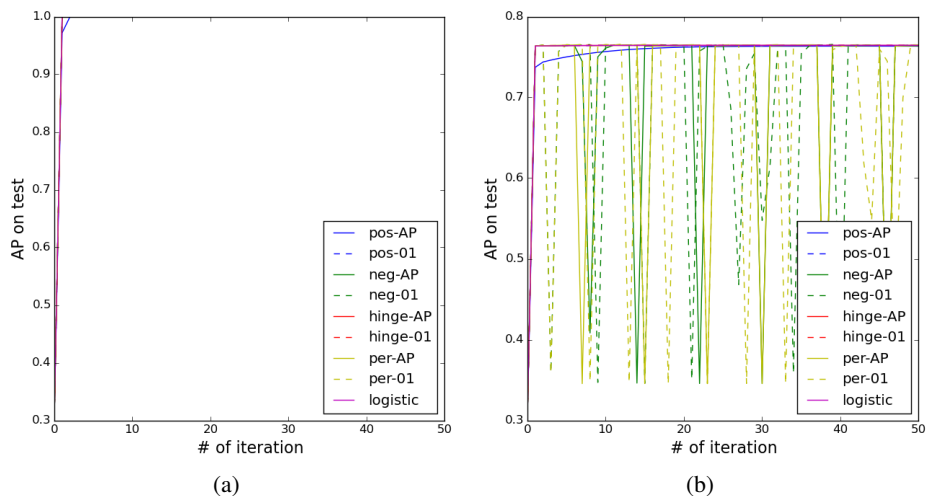
Figure 1: This figure shows the results on linear synthetic data. We illustrate average precision over the number of iterations on the test set without noise in (a) and with 20% noise in (b).

## 3 EXPERIMENTS ON LINEAR SYNTHETIC DATA

To test the linear case, we generated two different datasets, one of which is linearly separable while the other one is not. We randomly generated 20,000 data points by sampling from a 10 dimensional standard Gaussian distribution. The data points with a sum of numbers in all dimensions being larger than 0 are assigned to the positive class while those having a negative sum are classified as the negative objects. We then divide the whole dataset into training set and test set of 10,000 elements each. To produce the non-linearly separable dataset, we randomly flip 20% of the binary labels. We select $\phi(x, w) = w^\intercal x$ in this linear setting. The results are depicted in Fig. 1.

In the noiseless linear case we observe the negative update to achieve a slightly better performance than the positive update. The perceptron method also performs well. We think this is the reason why McAllester et al. (2010) report the negative update to perform better. Note that Cheng et al. (2009) also reported good performance for the perceptron method on the TIMIT dataset, the same one used in McAllester et al. (2010).

Negative and perceptron updates perform similarly on the noisy and not linearly separable dataset. They also do not perform well on our nonlinear datasets shown in the main paper.

## REFERENCES

Cheng, C.-C., Sha, F., and Saul, L. K. Matrix updates for perceptron training of continuous density hidden markov models. In *Proc. ICML*, 2009.

McAllester, D. A., Keshet, J., and Hazan, T. Direct loss minimization for structured prediction. In *Proc. NIPS*, 2010.