# Machine learning theory

## Ranking

Hamid Beigy

Sharif university of technology

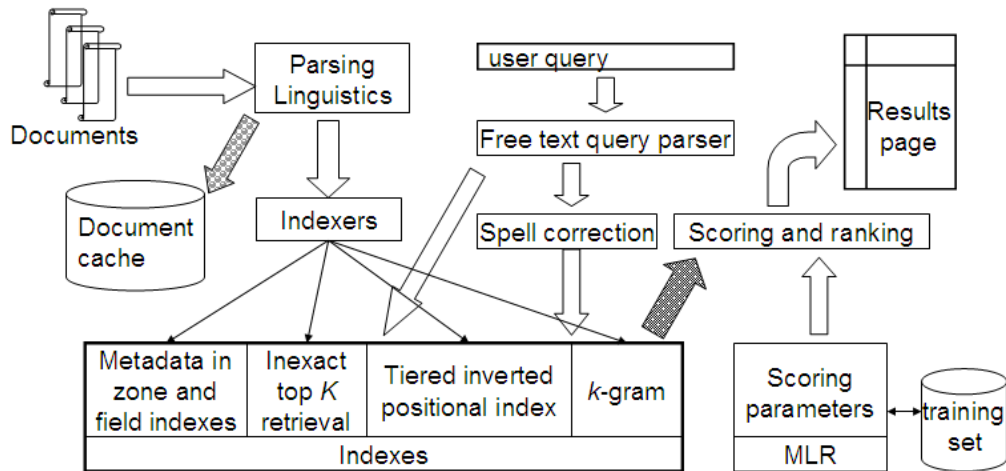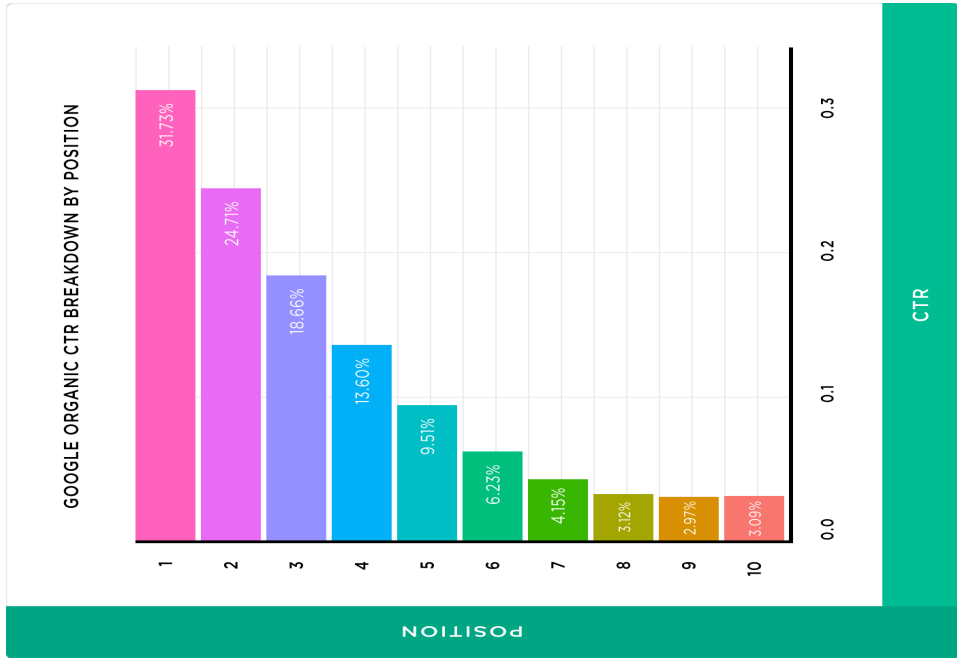May 19, 2023

# Introduction

A complete search engine

1. The first rank has average click rate of 31.7%.
2. Only 0.78% of Google searchers clicked from the second page.

1. The learning to rank problem is how to learn an ordering.
2. Application in very large datasets
   - search engines,
   - information retrieval
   - fraud detection
   - movie recommendation

---

**Motivation for ranking**

The main motivation for ranking over classification in the binary case is the limitation of resources.

1. it may be impractical or even impossible to display or process all items labeled as relevant by a classifier.
2. we need to show more relevant ones or prioritize them.

---

1. In applications such as search engines, ranking is more desirable than classification.

2. **Problem:** Can we learn to predict ranking accurately?

3. Ranking scenarios
   - score-based setting
   - preference-based setting

# Score-based setting

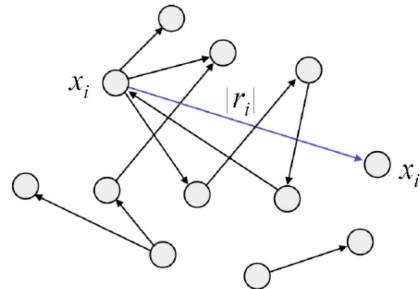General supervised learning problem of ranking,

- the learner receives labeled sample of pairwise preferences,
- the learner outputs a scoring function $h : \mathcal{X} \mapsto \mathbb{R}$.

**Drawbacks**

- $h$ induces a linear ordering for full set $\mathcal{X}$
- does not match a query-based scenario.

**Advantages**

- efficient algorithms
- good theory,
- VC bounds,
- margin bounds,
- stability bounds

1. The score-based setting is defined as
   - $\mathcal{X}$ is input space.
   - $\mathcal{D}$ is unknown distribution over $\mathcal{X} \times \mathcal{X}$.
   - $f : \mathcal{X} \times \mathcal{X} \mapsto \{-1, 0, +1\}$ is target labeling function or preference function, where

   $$f(\mathbf{x}, \mathbf{x}') = \begin{cases} -1 & \text{if } \mathbf{x}' \prec_{pref} \mathbf{x} \\ 0 & \text{if } \mathbf{x}' =_{pref} \mathbf{x} \\ +1 & \text{if } \mathbf{x} \prec_{pref} \mathbf{x}' \end{cases}$$

2. No assumption is made about the **transitivity** of the order induced by $f$.

   $$f(\mathbf{x}, \mathbf{x}') = +1 \quad \text{and} \quad f(\mathbf{x}', \mathbf{x}'') = +1 \quad \text{and} \quad f(\mathbf{x}'', \mathbf{x}) = +1$$

3. No assumption is made about the **antisymmetry** of the order induced

   $$f(\mathbf{x}, \mathbf{x}') = +1 \quad \text{and} \quad f(\mathbf{x}', \mathbf{x}) = +1 \quad \text{and} \quad \mathbf{x} \neq \mathbf{x}'$$

**Definition (Learning to rank (score-based setting))**

1. Learner receives $S = \{(\mathbf{x}_1, \mathbf{x}_1', y_1), \ldots, (\mathbf{x}_m, \mathbf{x}_m', y_m)\} \in (\mathcal{X} \times \mathcal{X} \mapsto \{-1, 0, +1\})^m$, where $(\mathbf{x}_i, \mathbf{x}_i') \sim \mathcal{D}$ and $y_i = f(\mathbf{x}_i, \mathbf{x}_i')$.

2. Given a hypothesis set $H = \{h : \mathcal{X} \mapsto \mathbb{R}\}$, ranking problem consists of selecting a hypothesis $h \in H$ with small expected pairwise misranking or generalization error $\mathbf{R}(h)$ with respect to the target $f$

$$\mathbf{R}(h) = \mathbb{P}_{(\mathbf{x}, \mathbf{x}') \sim \mathcal{D}} \left[ (f(\mathbf{x}, \mathbf{x}') \neq 0) \wedge (f(\mathbf{x}, \mathbf{x}')(h(\mathbf{x}) - h(\mathbf{x}')) \leq 0) \right]$$

3. The empirical pairwise misranking or empirical error of $h$ is defined by

$$\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^{m} \mathbb{I} \left[ (y_i \neq 0) \wedge (y_i(h(\mathbf{x}_i) - h(\mathbf{x}_i')) \leq 0) \right]$$

1. A simple approach is to project instances into a vector **w**

2. Let to define the ranking function as

$$h((\mathbf{x}_1, \ldots, \mathbf{x}_m)) = (\langle \mathbf{w}, \mathbf{x}_1 \rangle, \ldots, \langle \mathbf{w}, \mathbf{x}_m \rangle)$$

3. Then use the distance of the point to classifier $\langle \mathbf{w}, \mathbf{x} \rangle$ as the score of **x**.

4. We assume that $y_i \neq 0$, then the empirical error is defined as

$$\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\left[ (y_i(h(\mathbf{x}_i) - h(\mathbf{x}_i')) \leq 0) \right]$$

5. if we define $h(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle$, we have

$$\hat{\mathbf{R}}(h) = \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\left[ (y_i \langle \mathbf{w}, (\mathbf{x}_i - \mathbf{x}_i') \rangle \leq 0) \right]$$

6. Then, we can use the following ERM algorithm to rank items.

$$\mathbf{w} = \arg\min_{w'} \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\left[ (y_i \langle \mathbf{w}', (\mathbf{x}_i - \mathbf{x}_i') \rangle \leq 0) \right]$$
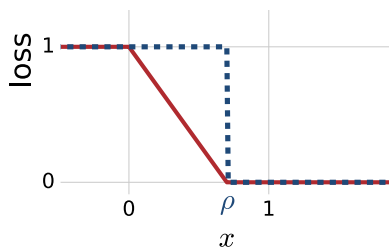
1. Assume that labels are chosen from $\{-1, +1\}$.
2. **Homework:** Generalize the result to the label set $\{-1, 0, +1\}$.
3. Same as classification, for any $\rho > 0$, empirical margin loss of a hypothesis $h$ for pairwise ranking is

$$\hat{\mathbf{R}}_\rho(h) = \frac{1}{m} \sum_{i=1}^{m} \Phi_\rho(y_i(h(\mathbf{x}_i') - h(\mathbf{x}_i)))$$

where

$$\Phi_\rho(u) = \begin{cases} 1 & \text{if } u \leq 0 \\ 1 - \dfrac{u}{\rho} & \text{if } 0 \leq u \leq \rho \\ 0 & \text{if } \rho \geq u \end{cases}$$
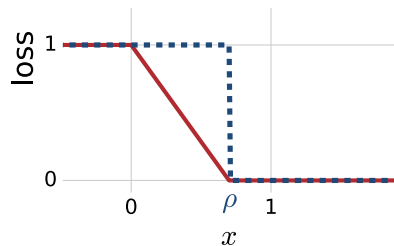
4. The parameter $\rho > 0$ can be interpreted as the confidence margin demanded from a hypothesis $h$.

The upper bound of empirical margin loss of a hypothesis $h$ is

$$\hat{\mathbf{R}}_\rho(h) \leq \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\left[y_i(h(\mathbf{x}_i') - h(\mathbf{x}_i)) \leq \rho\right]$$



Let

1. $\mathcal{D}_1$ be marginal distribution of the first element of pairs $\mathcal{X} \times \mathcal{X}$ derived from $\mathcal{D}$,
2. $\mathcal{D}_2$ be marginal distribution of the second element of pairs $\mathcal{X} \times \mathcal{X}$ derived from $\mathcal{D}$,
3. $S_1 = \{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m)\}$ and $\mathcal{R}_m^{\mathcal{D}_1}(H)$ be the Rademacher complexity of $H$ with respect to $\mathcal{D}_1$,
4. $S_2 = \{(\mathbf{x}_1', y_1), \ldots, (\mathbf{x}_m', y_m)\}$ and $\mathcal{R}_m^{\mathcal{D}_2}(H)$ be the Rademacher complexity of $H$ with respect to $\mathcal{D}_2$,

1. We also have $\mathcal{R}_m^{\mathcal{D}_1}(H) = \mathbb{E}\left[\hat{\mathcal{R}}_{S_1}(H)\right]$ and $\mathcal{R}_m^{\mathcal{D}_2}(H) = \mathbb{E}\left[\hat{\mathcal{R}}_{S_2}(H)\right]$.

2. If $\mathcal{D}$ is symmetric, then $\mathcal{R}_m^{\mathcal{D}_1}(H) = \mathcal{R}_m^{\mathcal{D}_2}(H)$.

---

**Theorem (Margin bound for ranking)**

*Let $H$ be a set of real-valued functions. Fix $\rho > 0$, then, for any $\delta > 0$, with probability at least $(1 - \delta)$ over the choice of a sample $S$ of size $m$, each of the following holds for all $h \in H$*

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}_\rho(h) + \frac{2}{\rho}\left(\mathcal{R}_m^{\mathcal{D}_1}(H) + \mathcal{R}_m^{\mathcal{D}_2}(H)\right) + \sqrt{\frac{\log(1/\delta)}{2m}}$$

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}_\rho(h) + \frac{2}{\rho}\left(\hat{\mathcal{R}}_{S_1}(H) + \hat{\mathcal{R}}_{S_2}(H)\right) + 3\sqrt{\frac{\log(2/\delta)}{2m}}$$

---

**Proof (Margin bound for ranking).**

1. Consider the family of functions $\tilde{H} = \{\Phi_\rho \circ h \mid f \in H\}$.

2. From margin-loss bounds we have

$$\mathbb{E}\left[\Phi_\rho(y[h(\mathbf{x}') - h(\mathbf{x})]\right] \leq \hat{\mathbf{R}}_\rho(h) + 2\mathcal{R}_m(\Phi_\rho \circ H) + \sqrt{\frac{\log(1/\delta)}{2m}}.$$

3. Since for all $u \in \mathbb{R}$, we have $\mathbb{I}[u \leq 0] \leq \Phi_\rho(u)$, then we have

$$\mathbf{R}(h) = \mathbb{E}\left[\mathbb{I}\left[y(h(\mathbf{x}') - h(\mathbf{x})) \leq 0\right]\right] \leq \mathbb{E}\left[\Phi_\rho(y[h(\mathbf{x}') - h(\mathbf{x})])\right]$$

4. Hence, we can write

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}_\rho(h) + 2\mathcal{R}_m(\Phi_\rho \circ H) + \sqrt{\frac{\log(1/\delta)}{2m}}.$$

5. Since $\Phi_\rho$ is $1/\rho - Lipschitz$, by Talagrand's lemma $\mathcal{R}_m(\Phi_\rho \circ \tilde{H}) \leq \frac{1}{\rho}\mathcal{R}_m(H)$.

**Proof (Margin bound for ranking)(cont.).**

6. Here, $\mathcal{R}_m(H)$ can be upper bounded as

$$
\begin{aligned}
\mathcal{R}_m(H) &= \frac{1}{m} \mathop{\mathbb{E}}_{S,\sigma} \left[ \sup_{h \in H} \sum_{i=1}^m \sigma_i y_i (h(\mathbf{x}_i') - h(\mathbf{x}_i)) \right] \\
&= \frac{1}{m} \mathop{\mathbb{E}}_{S,\sigma} \left[ \sup_{h \in H} \sum_{i=1}^m \sigma_i (h(\mathbf{x}_i') - h(\mathbf{x}_i)) \right] \quad \textcolor{magenta}{\sigma_i y_i \text{ and } \sigma_i : \text{same distribution}} \\
&\leq \frac{1}{m} \mathop{\mathbb{E}}_{S,\sigma} \left[ \sup_{h \in H} \sum_{i=1}^m \sigma_i h(\mathbf{x}_i') + \sup_{h \in H} \sum_{i=1}^m \sigma_i h(\mathbf{x}_i) \right] \quad \textcolor{magenta}{\text{by sub-additivity of sup}} \\
&\leq \mathop{\mathbb{E}}_{S} \left[ \hat{\mathcal{R}}_{S_1}(H) + \hat{\mathcal{R}}_{S_2}(H) \right] \quad \textcolor{magenta}{\text{definition of } S_1 \text{ and } S_2} \\
&\leq \mathcal{R}_m^{\mathcal{D}_1}(H) + \mathcal{R}_m^{\mathcal{D}_2}(H).
\end{aligned}
$$

7. The second inequality, can be derived in the same way.

$\square$

These bounds can be generalized to hold uniformly for any $\rho > 0$ at cost of an additional term $\sqrt{(\log \log_2(2/\rho))/m}$.

**Corollary (Margin bounds for ranking with kernel-based hypotheses)**

*Let $K : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ be a PDS kernel with $r = \sup_{\mathbf{x} \in \mathcal{X}} K(\mathbf{x}, \mathbf{x})$. Let also $\Phi : \mathcal{X} \mapsto \mathbb{H}$ be a feature mapping associated to $K$ and let $H = \left\{ \mathbf{x} \mapsto \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle \mid \|\mathbf{w}\|_{\mathbb{H}} \leq \Lambda \right\}$ for some $\Lambda \geq 0$. Fix $\rho > 0$. Then, for any $\delta > 0$, the following pairwise margin bound holds with probability at least $(1 - \delta)$ for any $h \in H$:*

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}_\rho(h) + 4\sqrt{\frac{r^2\Lambda^2/\rho^2}{m}} + \sqrt{\frac{\log(1/\delta)}{2m}}$$

1. This bound can be generalized to hold uniformly for any $\rho > 0$ at cost of an additional term $\sqrt{(\log \log_2(2/\rho))/m}$.
2. This bound suggests that a small generalization error can be achieved
   - when $\frac{\rho}{r}$ is large (small second term),
   - while the empirical margin loss is relatively small (first term).

From the generalization bound for SVM, Corollary Margin bounds for ranking with kernel-based hypotheses can be expressed as

---

**Corollary (Margin bounds for ranking with SVM)**

*Let $K : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ be a PDS kernel with $r = \sup_{\mathbf{x} \in \mathcal{X}} K(\mathbf{x}, \mathbf{x})$. Let also $\Phi : \mathcal{X} \mapsto \mathbb{H}$ be a feature mapping associated to $K$ and let $H = \left\{ \mathbf{x} \mapsto \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle \mid \|\mathbf{w}\|_{\mathbb{H}} \leq \Lambda \right\}$ for some $\Lambda \geq 0$. Then, for any $\delta > 0$, the following pairwise margin bound holds with probability at least $(1 - \delta)$ for any $h \in H$:*

$$\mathbf{R}(h) \leq \frac{1}{m} \sum_{i=1}^{m} \xi_i + 4\sqrt{\frac{r^2 \Lambda^2}{m}} + \sqrt{\frac{\log(1/\delta)}{2m}}$$

*where $\xi = \max\left(1 - y_i \left[\Phi(\mathbf{x}'_i) - \Phi(\mathbf{x}_i)\right], 0\right)$*

---

1. Margin bounds for ranking with SVM

$$\mathbf{R}(h) \leq \frac{1}{m} \sum_{i=1}^{m} \xi_i + 4\sqrt{\frac{r^2 \Lambda^2}{m}} + \sqrt{\frac{\log(1/\delta)}{2m}}$$

2. Minimizing the right-hand side of this inequality is
   minimizing an objective function with a term corresponding to the sum of the slack variables $\xi_i$,
   and another one minimizing $\|\mathbf{w}\|$ or equivalently $\|\mathbf{w}\|^2$.

3. This optimization problem can thus be formulated as

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^{m} \xi_i$$

$$\text{subject to } y_i \left[ \langle \mathbf{w}, \left( \Phi(\mathbf{x}_i') - \Phi(\mathbf{x}_i) \right) \rangle \right] \geq 1 - \xi_i$$

$$\xi_i \geq 0 \quad \forall 1 \leq i \leq m.$$

1. This optimization problem coincides exactly with the primal optimization problem of SVMs, with a feature mapping

$$\Psi : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{H}$$

   defined by

$$\Psi(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) - \Phi(\mathbf{x}')$$

   for all

$$(\mathbf{x}, \mathbf{x}') \in \mathcal{X} \times \mathcal{X}$$

   and with a hypothesis set of functions of the form

$$(\mathbf{x}, \mathbf{x}') \mapsto \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{x}') \rangle .$$

2. Clearly, all the properties already presented for SVMs apply in this instance.

3. In particular, the algorithm can benefit from the use of PDS kernels.

4. This can be used with kernels

$$
\begin{aligned}
K'((\mathbf{x}_i, \mathbf{x}_i'), (\mathbf{x}_j, \mathbf{x}_j')) &= \langle \Psi(\mathbf{x}_i, \mathbf{x}_i'), \Psi(\mathbf{x}_j, \mathbf{x}_j') \rangle \\
&= K(\mathbf{x}_i, \mathbf{x}_j) + K(\mathbf{x}_i', \mathbf{x}_j') - K(\mathbf{x}_i', \mathbf{x}_j) - K(\mathbf{x}_i, \mathbf{x}_j').
\end{aligned}
$$

**Boosting for ranking**

## Boosting for ranking

- Use weak ranking algorithm and create stronger ranking algorithm:
- Ensemble method: combine base rankers returned by weak ranking algorithm
- Finding simple relatively accurate base rankers often not hard.
- How should base rankers be combined?
- Let $H$ defined as

$$H = \{h : \mathcal{X} \mapsto \{0, 1\}\}$$

  where $H$ is the hypothesis set from which the base rankers are selected.

- For any $s \in \{-1, 0, +1\}$, we define

$$\epsilon_t^s = \sum_{i=1}^{m} D_t(i) \, \mathbb{I}\left[y_i(h_t(\mathbf{x}_i') - h_t(\mathbf{x}_i)) = s\right] = \mathop{\mathbb{E}}_{i \sim D_t}\left[\mathbb{I}\left[y_i(h_t(\mathbf{x}_i') - h_t(\mathbf{x}_i)) = s\right]\right]$$

- Hence, we have

$$\epsilon_t^+ + \epsilon_t^- + \epsilon_t^0 = 1$$

- We assume that $y_i \neq 0$.
- **Homework:** Show that the derivation of the algorithm.

**RankBoost Algorithm**

1: **function** RANKBOOST($S$, $H$, $T$)
2:    **for** $i \leftarrow 1$ to $m$ **do**
3:        $D_1(i) \leftarrow \dfrac{1}{m}$
4:    **end for**
5:    **for** $t \leftarrow 1$ to $T$ **do**
6:        Let $h_t = \arg\min_{h \in H} \left( \epsilon_t^- - \epsilon_t^+ \right)$

        $\triangleright \epsilon^-$ : pairwise ranking error
        $\triangleright \epsilon^+$ : pairwise ranking accuracy

7:        $\alpha_t \leftarrow \dfrac{1}{2} \log \dfrac{\epsilon_t^+}{\epsilon_t^-}$
8:        $Z_t \leftarrow \epsilon_t^0 + 2\sqrt{\epsilon_t^+ \epsilon_t^-}$
9:        **for** $i \leftarrow 1$ to $m$ **do**
10:           $D_{t+1}(i) \leftarrow \dfrac{D_t(i)\exp\left[-\alpha_t y_i \left( h_t(\mathbf{x}_i') - h_t(\mathbf{x}_i) \right)\right]}{Z_t}$
11:       **end for**
12:    **end for**
13:    **return** $f \triangleq \sum_{t=1}^{T} \alpha_t h_t$
14: **end function**

**Theorem (Bound on the empirical error of RankBoost)**

*The empirical error of the hypothesis $H = \{h : \mathcal{X} \mapsto \{0,1\}\}$ returned by RankBoost verifies:*

$$\hat{\mathbf{R}}(h) \leq \exp\left[-2 \sum_{t=1}^{T} \left(\frac{\epsilon_t^+ - \epsilon_t^-}{2}\right)^2\right]$$
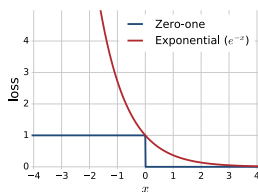
*Furthermore, if there exists $\gamma$ such that for all $1 \leq t \leq T$, condition $0 \leq \gamma \leq \dfrac{\epsilon_t^+ - \epsilon_t^-}{2}$, then*

$$\hat{\mathbf{R}}(h) \leq \exp\left[-2\gamma^2 T\right].$$

**Proof of (Bound on the empirical error of RankBoost).**

1. The empirical error equals to $\hat{R}(h) = \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\left[y_i(f(\mathbf{x}_i') - f(\mathbf{x}_i)) \leq 0\right]$.

2. On the other hand, for all $u \in \mathbb{R}$, we have $\mathbb{I}\left[u \leq 0\right] \leq \exp(-u)$.



3. Hence, we can write

$$\hat{R}(h) = \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\left[y_i(f(\mathbf{x}_i') - f(\mathbf{x}_i)) \leq 0\right]$$

$$\leq \frac{1}{m} \sum_{i=1}^{m} \exp\left[-y_i(f(\mathbf{x}_i') - f(\mathbf{x}_i))\right]$$

$$\leq \frac{1}{m} \sum_{i=1}^{m} \left[m \prod_{t=1}^{T} Z_t\right] D_{t+1}(i) = \prod_{t=1}^{T} Z_t.$$

**Proof of (Bound on the empirical error of RankBoost) (cont.).**

4. From definition of
$$Z_t = \sum_{i=1}^{m} D_t(i) exp\left[-y_i(h_t(\mathbf{x}_i') - h_t(\mathbf{x}_i))\right]$$
,

5. By grouping together the indices $i$ for which $y_i(h_t(\mathbf{x}_i') - h_t(\mathbf{x}_i))$ take values in $-1$, $0$, or $+1$, $Z_t$ can be written as

$$Z_t = \epsilon_t^+ e^{-\alpha_t} + \epsilon_t^- e^{+\alpha_t} + \epsilon_t^0$$

$$= \epsilon_t^+ \sqrt{\frac{\epsilon_t^-}{\epsilon_t^+}} + \epsilon_t^- \sqrt{\frac{\epsilon_t^+}{\epsilon_t^-}} + \epsilon_t^0$$

$$= 2\sqrt{\epsilon_t^+ \epsilon_t^-} + \epsilon_t^0$$

6. Since, $\epsilon_t^+ = 1 - \epsilon_t^- - \epsilon_t^0$, we have

$$4\epsilon_t^+ \epsilon_t^- = \left(\epsilon_t^+ + \epsilon_t^-\right)^2 - \left(\epsilon_t^+ - \epsilon_t^-\right)^2 = \left(1 - \epsilon_t^0\right)^2 - \left(\epsilon_t^+ - \epsilon_t^-\right)^2$$

**Proof of (Bound on the empirical error of RankBoost) (cont.).**

7. Thus, assuming that $\epsilon_t^0 < 1$, $Z_t$ can be upper bounded as

$$Z_t = \sqrt{(1 - \epsilon_t^0)^2 - \left(\epsilon_t^+ - \epsilon_t^-\right)^2} + \epsilon_t^0 = \left(1 - \epsilon_t^0\right)\sqrt{1 - \frac{\left(\epsilon_t^+ - \epsilon_t^-\right)^2}{\left(1 - \epsilon_t^0\right)^2}} + \epsilon_t^0$$

$$\leq \left(1 - \epsilon_t^0\right)\exp\left(-\frac{\left(\epsilon_t^+ - \epsilon_t^-\right)^2}{2\left(1 - \epsilon_t^0\right)^2}\right) + \epsilon_t^0 \qquad \text{By using inequality } 1 - x \leq e^{-x}$$

$$\leq \exp\left(-\frac{\left(\epsilon_t^+ - \epsilon_t^-\right)^2}{2}\right) \qquad\qquad \text{exp is concave and } 0 < \left(1 - \epsilon_t^0\right) \leq 1$$

$$\leq \exp\left(-2\left[\frac{\left(\epsilon_t^+ - \epsilon_t^-\right)}{2}\right]^2\right)$$

8. By setting $0 \leq \gamma \leq \dfrac{\epsilon_t^+ - \epsilon_t^-}{2}$, we obtain $\hat{\mathbf{R}}(h) \leq \exp\left[-2\gamma^2 T\right]$.

$\square$

1. Assume that the pairwise labels are in $\{-1, +1\}$.
2. We showed that $\hat{\mathcal{R}}_S(conv(H)) = \hat{\mathcal{R}}_S(H)$.

---

**Corollary (Margin bound for ensemble methods in ranking)**

*Let $H$ be a set of real-valued functions. Fix $\rho > 0$; then, for any $\delta > 0$, with probability at least $(1 - \delta)$ over the choice of a sample $S$ of size $m$, each of the following ranking guarantees holds for all $h \in conv(H)$*

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}_\rho(h) + \frac{2}{\rho}\left(\mathcal{R}_m^{\mathcal{D}_1}(H) + \mathcal{R}_m^{\mathcal{D}_2}(H)\right) + \sqrt{\frac{\log(1/\delta)}{2m}}$$

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}_\rho(h) + \frac{2}{\rho}\left(\hat{\mathcal{R}}_{S_1}(H) + \hat{\mathcal{R}}_{S_2}(H)\right) + 3\sqrt{\frac{\log(2/\delta)}{2m}}$$

---

3. These bounds apply to $h/\|\alpha\|_1$, where $h$ and $h/\|\alpha\|_1$ induce the same ordering.
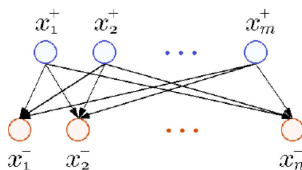4. Then, or any $\delta > 0$, the following holds with probability at least $(1 - \delta)$

$$\mathbf{R}(h) \leq \hat{\mathbf{R}}_\rho(h/\|\alpha\|_1) + \frac{2}{\rho}\left(\mathcal{R}_m^{\mathcal{D}_1}(H) + \mathcal{R}_m^{\mathcal{D}_2}(H)\right) + \sqrt{\frac{\log(1/\delta)}{2m}}$$

5. Note that $T$ does not appear in this bound.

**Bipartite ranking**

1. Bipartite ranking problem is an important ranking scenario within score-based setting.
2. In this scenario, the set of points $\mathcal{X}$ is partitioned into
   - the class of positive points $\mathcal{X}_+$
   - the class of negative points $\mathcal{X}_-$
3. In this setting, positive points must rank higher than negative ones and the learner receives
   - a sample $S_+ = (\mathbf{x}_1', \ldots, \mathbf{x}_m')$ drawn i.i.d. according to some distribution $\mathcal{D}_+$ over $\mathcal{X}_+$ ,
   - a sample $S_- = (\mathbf{x}_1, \ldots, \mathbf{x}_n)$ drawn i.i.d. according to some distribution $\mathcal{D}_-$ over $\mathcal{X}_-$ .

1. The learning problem consists of selecting a hypothesis $h \in H$ with small expected bipartite misranking or generalization error $\mathbf{R}(h)$ :

$$\mathbf{R}(h) = \mathbb{P}_{\substack{\mathbf{x}' \sim \mathcal{D}_+ \\ \mathbf{x} \sim \mathcal{D}_-}} \left[ h(\mathbf{x}') < h(\mathbf{x}) \right]$$

2. The empirical pairwise mis-ranking or empirical error of $h$ is

$$\hat{\mathbf{R}}_{S_+, S_-}(h) = \frac{1}{mn} \sum_{i=1}^{m} \sum_{j=1}^{n} \mathbb{I} \left[ h(\mathbf{x}'_i) < h(\mathbf{x}_j) \right]$$

3. The learning algorithm must typically deal with $mn$ pairs.

1. A key property of RankBoost leading to an efficient algorithm for bipartite ranking is exponential form of its objective function.
2. The objective function can be decomposed into the product of two functions,
   - one depends on only the positive points.
   - one depends on only the negative points.
3. Similarly,

$$
\begin{aligned}
D_1(i,j) &= \frac{1}{mn} \\
&= D_1^+(i) D_1^-(j) \\
&= \frac{1}{m} \times \frac{1}{n}
\end{aligned}
$$

4. Similarly,

$$
\begin{aligned}
D_{t+1}(i,j) &= \frac{D_t(i,j) \exp\left(-\alpha_t \left[h_t(\mathbf{x}_i') - h_t(\mathbf{x}_j)\right]\right)}{Z_t} \\
&= \frac{D_t^+(i) \exp\left(-\alpha_t h_t(\mathbf{x}_i')\right)}{Z_t^+} \times \frac{D_t^-(j) \exp\left(\alpha_t h_t(\mathbf{x}_j)\right)}{Z_t^-}
\end{aligned}
$$

1. The pairwise misranking of a hypothesis $h$

$$
\begin{aligned}
\left(\epsilon_t^- - \epsilon_t^+\right) &= \mathop{\mathbb{E}}_{(i,j)\sim D_t} \left[h(\mathbf{x}_i') - h(\mathbf{x}_j)\right] \\
&= \mathop{\mathbb{E}}_{i\sim D_t^+} \left[\mathop{\mathbb{E}}_{j\sim D_t^-} \left[h(\mathbf{x}_i') - h(\mathbf{x}_j)\right]\right] \\
&= \mathop{\mathbb{E}}_{j\sim D_t^+} \left[h(\mathbf{x}_j')\right] - \mathop{\mathbb{E}}_{i\sim D_t^-} \left[h(\mathbf{x}_i)\right]
\end{aligned}
$$

2. The time and space complexity of BipartiteRankBoost is $O(m+n)$.

**BipartiteRankBoost Algorithm**

1: **function** BIPARTITERANKBOOST($S$, $H$, $T$)
2: $\quad D_1^+(i) \leftarrow \dfrac{1}{m} \quad \forall i \in 1, 2, \ldots, m$
3: $\quad D_1^-(j) \leftarrow \dfrac{1}{n} \quad \forall j \in 1, 2, \ldots, n$
4: $\quad$ **for** $t \leftarrow 1$ to $T$ **do**
5: $\quad\quad$ Let $h_t = \arg\min_{h \in H} \left( \epsilon_t^- - \epsilon_t^+ \right)$
6: $\quad\quad \alpha_t \leftarrow \dfrac{1}{2} \log \dfrac{\epsilon_t^+}{\epsilon_t^-}$
7: $\quad\quad Z_t^+ \leftarrow 1 - \epsilon_t^+ + \sqrt{\epsilon_t^+ \epsilon_t^-}$
8: $\quad\quad$ **for** $i \leftarrow 1$ to $m$ **do**
9: $\quad\quad\quad D_{t+1}^+(i) \leftarrow \dfrac{D_t^+(i) \exp\left[-\alpha_t h_t(\mathbf{x}_i')\right]}{Z_t^+}$
10: $\quad\quad$ **end for**
11: $\quad\quad Z_t^- \leftarrow 1 - \epsilon_t^- + \sqrt{\epsilon_t^+ \epsilon_t^-}$
12: $\quad\quad$ **for** $j \leftarrow 1$ to $n$ **do**
13: $\quad\quad\quad D_{t+1}^-(j) \leftarrow \dfrac{D_t^-(j) \exp\left[\alpha_t h_t(\mathbf{x}_j)\right]}{Z_t^-}$
14: $\quad\quad$ **end for**
15: $\quad$ **end for**
16: $\quad$ **return** $f \triangleq \sum_{t=1}^T \alpha_t h_t$
17: **end function**

1. The objective function of RankBoost can be expressed as

$$F_{RankBoost}(\alpha) = \sum_{j=1}^{m} \sum_{i=1}^{n} \exp\left(-\left[f(x_j') - f(x_j)\right]\right)$$

$$= \left(\sum_{i=1}^{m} \exp\left(-\sum_{t=1}^{T} \alpha_t h_t(x_i')\right)\right) \left(\sum_{j=1}^{n} \exp\left(\sum_{t=1}^{T} \alpha_t h_t(x_j)\right)\right)$$

$$= F_{+}(\alpha) F_{-}(\alpha)$$

where $F_{+}(\alpha)$ denotes function defined by the sum over positive points and $F_{-}(\alpha)$ function defined over negative points.
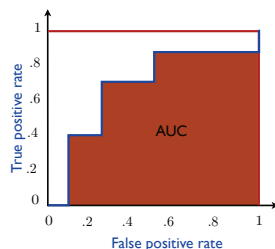
2. The objective function of AdaBoost can be expressed as

$$F_{AdaBoost}(\alpha) = \sum_{j=1}^{m} \exp\left(-y_j' f(x_j')\right) + \sum_{i=1}^{n} \exp\left(-y_i f(x_i)\right)$$

$$= \sum_{i=1}^{m} \exp\left(-\sum_{t=1}^{T} \alpha_t h_t(x_i')\right) + \sum_{j=1}^{n} \exp\left(\sum_{t=1}^{T} \alpha_t h_t(x_j)\right)$$

$$= F_{+}(\alpha) + F_{-}(\alpha)$$

1. Performance of a bipartite ranking algorithm is reported in terms of area ROC curve, or AUC.

2. Let $U$ be a test sample used for evaluating the performance of $h$

   - $m$ positive points $\mathbf{z}_1', \ldots, \mathbf{z}_m'$

   - $n$ negative points $\mathbf{z}_1, \ldots, \mathbf{z}_n$

   - $AUC(h, u)$ equals to

   $$AUC(h, U) = \frac{1}{mn} \sum_{i=1}^{m} \sum_{j=1}^{n} \mathbb{I}\left[h(\mathbf{z}_i') \geq h(\mathbf{z}_j)\right]$$

   $$= \mathop{\mathbb{P}}_{\substack{\mathbf{z} \sim D_U^- \\ \mathbf{z}' \sim D_U^+}} \left[h(\mathbf{z}') \geq h(\mathbf{z})\right]$$



3. The average pairwise misranking of $h$ over $U$ denoted by $\hat{\mathbf{R}}(h, U)$

$$\hat{\mathbf{R}}(h, U) = 1 - AUC(h, U).$$

4. AUC can be computed in time of $O(m + n)$ from a sorted array $h(\mathbf{z}_i')$ and $h(\mathbf{z}_j)$.

5. Homework: Design an algorithm for computing AUC in time of $O(m + n)$.

# Preference-based setting

1. Assume that you receive a list $X \subseteq \mathcal{X}$ as a result of a query $q$.

2. The goal is to rank items in list $X$ not all items in $\mathcal{X}$.

3. The advantage of preference-based setting over score-based setting is:
   **The learning algorithm is not required to return a linear ordering of all points of $\mathcal{X}$, which may be impossible.**

4. The preference-based setting consists of two stages.
   - A sample of labeled pairs $S$ is used to learn a **preference function** $h : \mathcal{X} \times \mathcal{X} \mapsto [0, 1]$.
   - Given list $X \subseteq \mathcal{X}$, the preference function $h$ is used to determine a ranking of $X$.

5. How can $h$ be used to generate an accurate ranking?

6. The computational complexity of the second stage is also crucial.

7. We will measure the time complexity in terms of the number of calls to $h$.

1. Assume that a preference function $h$ is given.

2. $h$ is not assumed to be transitive.

3. We assume that $h$ is pairwise consistent, that is

$$h(u, v) + h(v, u) = 1, \qquad \forall u, v \in \mathcal{X}$$

4. Let $\mathcal{D}$ be an unknown distribution according to which pairs $(X, \sigma^*)$ are drawn, where
   - $X \subseteq \mathcal{X}$ is a query subset.
   - $\sigma^*$ is a target ranking.

5. The objective of a second-stage algorithm $A$ is using function $h$ to return an accurate ranking $A(X)$ for any query subset $X$.

6. The algorithm $A$ may be deterministic or randomized.

1. Loss function $\ell$ is used to measure disagreement between target ranking $\sigma^*$ and ranking $\sigma$ for set $X$ with $n \geq 1$ elements.

$$\ell(\sigma, \sigma^*) = \frac{2}{n(n-1)} \sum_{u \neq v} \mathbb{I}\left[\sigma(u) < \sigma(v)\right] \mathbb{I}\left[\sigma^*(v) < \sigma^*(u)\right]$$

2. Loss between target ranking $\sigma^*$ and ranking $h$ equals to

$$\ell(h, \sigma^*) = \frac{2}{n(n-1)} \sum_{u \neq v} h(u, v) \mathbb{I}\left[\sigma^*(v) < \sigma^*(u)\right]$$

- The expected loss for a deterministic algorithm $A$ is

$$\mathop{\mathbb{E}}_{(X,\sigma^*)\sim\mathcal{D}}[\ell(A(X),\sigma^*)].$$

- Regret of algorithm $A$ is the difference between its loss and loss of the best fixed global ranking.

$$Regret(A) = \mathop{\mathbb{E}}_{(X,\sigma^*)\sim\mathcal{D}}[\ell(A(X),\sigma^*)] - \min_{\sigma'} \mathop{\mathbb{E}}_{(X,\sigma^*)\sim\mathcal{D}}\left[\ell(\sigma'_{|X},\sigma^*)\right]$$

- Regret of the preference function is

$$Regret(h) = \mathop{\mathbb{E}}_{(X,\sigma^*)\sim\mathcal{D}}\left[\ell(h_{|X},\sigma^*)\right] - \min_{h'} \mathop{\mathbb{E}}_{(X,\sigma^*)\sim\mathcal{D}}\left[\ell(h'_{|X},\sigma^*)\right]$$

1. For sort by degree algorithm $A$, we can prove

$$Regret(A) \leq 2Regret(h)$$
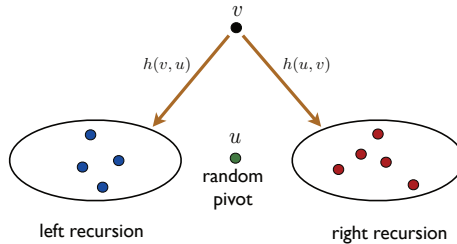
**Theorem (Lower bound for deterministic algorithms)**

*For any deterministic algorithm A, there is a bipartite distribution for which*

$$Regret(A) \geq 2Regret(h)$$

2. **Homework:** Prove the above theorem.

1. The second stage use a straightforward extension of the randomized QuickSort algorithm.



2. For randomized quick sort(RQS), we can prove

$$Regret(A_{RQS}) \leq Regret(h)$$

3. **Homework:** Prove the above bound.
4. **Homework:** Calculate the computation time of this algorithm.

**Extension to other loss functions**

1. All of the results just presented hold for a broader class of loss functions $L_w$ defined in terms of a weight function $w$.

$$L_w(\sigma, \sigma^*) = \frac{2}{n(n-1)} \sum_{u \neq v} w(\sigma^*(v) - \sigma^*(u)) \, \mathbb{I}\,[\sigma(u) < \sigma(v)] \, \mathbb{I}\,[\sigma^*(v) < \sigma^*(u)]$$

2. Function $w$ is assumed to satisfy the following three natural axioms:

   **Symmetry**  $w(i,j) = w(j,i)$ for all $i, j$.
   **Monotonicity**  $w(i,j) \leq w(i,k)$ if either $i < j < k$ or $i > j > k$.
   **Triangle inequality**  $w(i,j) \leq w(i,k) + w(k,j)$.

3. Using different functions $w$, the family of functions $L_w$ can cover several familiar and important losses.

# Summary

- We defined ranking problem.
- We extend this by using other loss functions defined in terms of a weight function.
- We can extend this by using other criteria have been introduced in information retrieval such as *NDCG*, *P@n*.

1. Sections 17.4 and 17.5 of Shai Shalev-Shwartz and Shai Ben-David (2014). *Understanding machine learning: From theory to algorithms*. Cambridge University Press.

2. Chapter 10 of Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar (2018). *Foundations of Machine Learning*. Second Edition. MIT Press.

3. The interested reader is referred to **Hang11** .

📄 Mohri, Mehryar, Afshin Rostamizadeh, and Ameet Talwalkar (2018). *Foundations of Machine Learning*. Second Edition. MIT Press.

📄 Shalev-Shwartz, Shai and Shai Ben-David (2014). *Understanding machine learning: From theory to algorithms*. Cambridge University Press.

**Questions?**