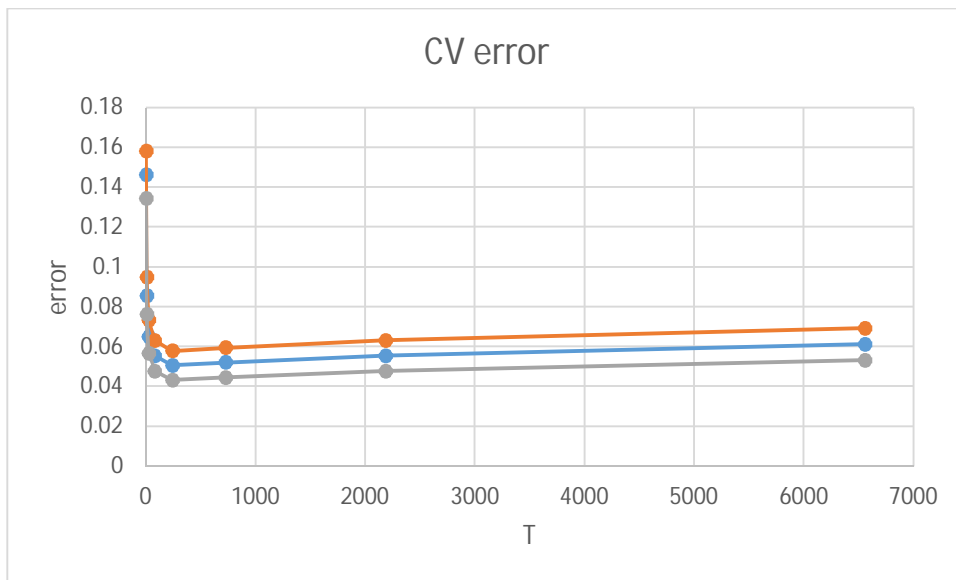


Paul Lintilhac

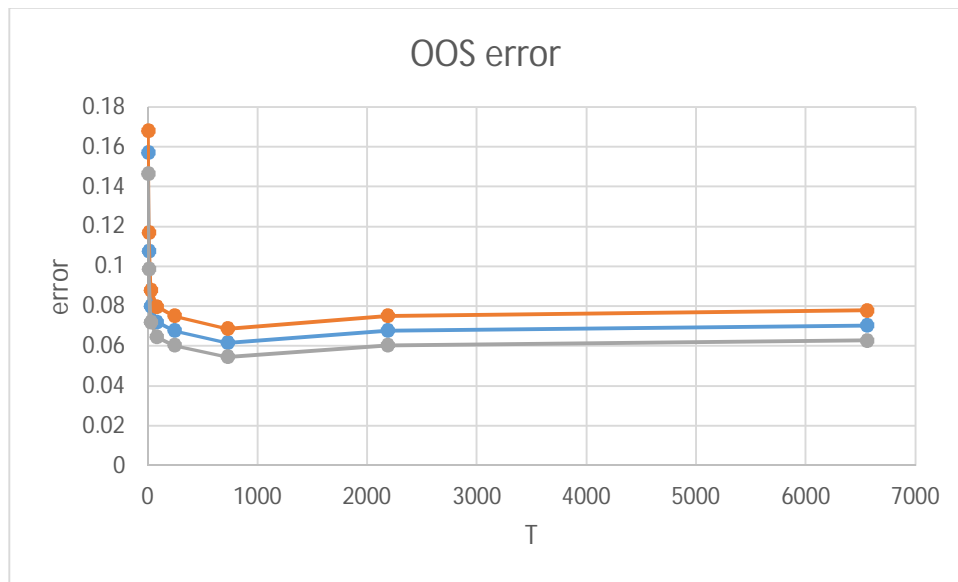
Homework 3 – Part A

Foundations Of Machine Learning

- 1) I did the problem without using any packages, and did a full implementation in both R and C++. The C++ is much faster to run of course, though the r code is easier to read. For the parameter T , I used powers of 3, i.e. $T_i = \{3^i: i \in \{1,2,3,4,5,6,7,8\}\}$. I used 4-fold cross-validation. Note that I calculated the error bars by first calculating the point-wise squared error: $\hat{\sigma}^2 = E[R(h)^2] - \hat{R}(h)^2 = \hat{R}(h) - \hat{R}(h)^2$, where I have used the fact that for a binary variable, $\hat{R}(h)^2 = \hat{R}(h)$. Next, since the observations are i.i.d., we can use the law of averages to compute the error over 3540/4 observations (for the cv error) by taking $\hat{\sigma}^2_{Tot} = 4 * \hat{\sigma}^2 / 3450$, or we can take the error over 1170 observations (for the OOS error) by taking $\hat{\sigma}^2_{Tot} = 4 * \hat{\sigma}^2 / 1151$.



Paul Lintilhac



Compared to the previous homework, I was able to achieve a lower error for both cross-validation and out-of-sample testing. The minimum error for cross-validation was 0.050435, achieved at $T=243$, while the minimum error for out-of-sample testing was 0.061686, achieved at $T=729$.

Paul Lintilhac

2)

a.

I will derive the modified boosting algorithm starting from the following assumptions:

- 1) Our empirical loss function is now $\hat{R}(h) = \mathbf{1}_{g \neq y < 0}$, i.e. if the guess was not correct but either the prediction or the response is 0 (not sure), then we don't count it towards the error.
- 2) We maintain the assumption that our final hypothesis is of the form $g = g_T = \sum_{t=1}^T \alpha_t h_t$.
- 3) We maintain the assumption that our distribution evolves as $D_{t+1}(i) = D_t(i) * \frac{e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$

Note that the equations determining α_t and Z_t are not yet known.

Using these assumptions, we can see that the argument for the upper bound on the empirical error is exactly the same, and therefore leads to the same equation:

$$\hat{R}(h) \leq \prod_{t=1}^T Z_t$$

Since Z_t is a normalization factor, we can again use the exact same identity

$$Z_t = \sum_{i=1}^m D_t(i) e^{-\alpha_t y_i h_t(x_i)}$$

However, this time we decompose the sum into 3 different sums:

$$\begin{aligned} Z_t &= \sum_{i: y_i h_t(x_i) = -1}^m D_t(i) e^{-\alpha_t y_i h_t(x_i)} + \sum_{i: y_i h_t(x_i) = 0}^m D_t(i) e^{-\alpha_t y_i h_t(x_i)} + \sum_{i: y_i h_t(x_i) = 1}^m D_t(i) e^{-\alpha_t y_i h_t(x_i)} \\ &= \epsilon_t^{-1} e^{\alpha_t} + \epsilon_t^0 + \epsilon_t^1 e^{-\alpha_t} \end{aligned}$$

Note that $\epsilon_t^0 = 1 - \epsilon_t^{-1} - \epsilon_t^1$, and thus we can re-write this as

$$Z_t = \epsilon_t^{-1} (e^{\alpha_t} - 1) + \epsilon_t^1 (e^{-\alpha_t} - 1)$$

Now, in order to minimize the empirical error, we can note that Z_t is convex and differentiable, and again minimize it with respect to α_t , which yields

$$\alpha_t = \frac{1}{2} \log \left(\frac{\epsilon_t^{-1}}{\epsilon_t^1} \right)$$

Now, substituting this back into the equation for the empirical error, we find that

$$\hat{R}(h) \leq \prod_{t=1}^T Z_t = \prod_{t=1}^T [\epsilon_t^{-1} (e^{\alpha_t} - 1) + \epsilon_t^1 (e^{-\alpha_t} - 1) + 1]$$

$$\begin{aligned}
&= \prod_{t=1}^T \left[\epsilon_t^{-1} \left(\sqrt{\frac{\epsilon_t^1}{\epsilon_t^{-1}}} - 1 \right) + \epsilon_t^1 \left(\sqrt{\frac{\epsilon_t^{-1}}{\epsilon_t^1}} - 1 \right) + 1 \right] \\
&= \prod_{t=1}^T \left[2\sqrt{\epsilon_t^1 \epsilon_t^{-1}} - \epsilon_t^1 - \epsilon_t^{-1} + 1 \right] \\
&= \prod_{t=1}^T \left[1 - \left(\sqrt{\epsilon_t^{-1}} - \sqrt{\epsilon_t^1} \right)^2 \right] \\
&\leq \prod_{t=1}^T e^{-\left(\sqrt{\epsilon_t^{-1}} - \sqrt{\epsilon_t^1} \right)^2} \quad (1)
\end{aligned}$$

(the above answers d)

Paul Lintilhac

b.

We can see from this definition that rather than in the binary example, where the exponent is positive as long as $|\epsilon_t^{-1} - \frac{1}{2}| > 0$, in this case the exponent is positive as long as $|\epsilon_t^1 - \epsilon_t^{-1}| > 0$, i.e. as long as our rate of error is less than our rate of accuracy (excluding any points where $y_i h_t(x_i) = 0$). Thus, our weak learning assumption is now that there exists an algorithm A, $\gamma > 0$, and a polynomial function $\text{poly}(\dots)$ such that for any $\delta > 0$ and any distribution D on X and for any target concept C, the following holds for any sample size $m \geq \text{poly}\left(\frac{1}{\delta}, n, \text{size}(C)\right)$:

$$P_{S \sim D^m}[R^{-1}(h) \leq R^1(h) + \gamma] \geq 1 - \delta$$

Where $R^1(h) = E[\mathbf{1}_{g*y=1}]$, and $R^{-1}(h) = E[\mathbf{1}_{g*y=-1}]$.

c.

The pseudocode for the algorithm is as follows:

For i=1 to m do:

$$D_1(i) = \frac{1}{m}$$

For t=1 to T do:

$h_t < -$ base classifier with small *relative* error: $\epsilon_t^{-1} - \epsilon_t^1$

$$\alpha_t = \frac{1}{2} \log\left(\frac{\epsilon_t^1}{\epsilon_t^{-1}}\right)$$

$$Z_t = 2\sqrt{\epsilon_t^1 \epsilon_t^{-1}} + \epsilon_t^0$$

For i in 1 to m do:

$$D_{t+1}(i) = D_t(i) e^{-\alpha_t * h_t(x_i) * y_i} / Z_t$$

$$g = \sum_{t=1}^T \alpha_t h_t$$

Paul Lintilhac

Part B:

1)

We start with the definition of $\Phi(\mathbf{x})$:

$$\Phi(\mathbf{x}) = \|\mathbf{x}_+\|_\alpha^2 = \left[\sum_{i=1}^N (x_i)_+^\alpha \right]^{\frac{2}{\alpha}}$$

$$\frac{\partial \Phi}{\partial x_i} = 2 \left[\sum_{j=1}^N (x_j)_+^\alpha \right]^{\frac{2}{\alpha}-1} (x_i)_+^{\alpha-1}$$

$$\frac{\partial^2 \Phi}{\partial x_i \partial x_j} = (4 - 2\alpha) \left[\sum_{k=1}^N (x_k)_+^\alpha \right]^{\frac{2}{\alpha}-2} (x_i)_+^{\alpha-1} (x_j)_+^{\alpha-1} + \left[(2\alpha - 2) \left[\sum_{k=1}^N (x_k)_+^\alpha \right]^{\frac{2}{\alpha}-1} (x_j)_+^{\alpha-2} \right]_{i=j}$$

where that the second term vanishes unless $i=j$.

Note that the lowest power of any individual $(x_j)_+$ is $\alpha - 2$. Since by assumption $\alpha > 2$, this means that the above mixed derivative (which can be used to compute the Hessian) is well-defined even when any particular $(x_j)_+ = 0$. However, observe what happens when all of the $(x_j)_+ = 0$, i.e. when $\mathbf{x} \in \mathbf{B}$. In this case,

$$\sum_{i=1}^N (x_i)_+^\alpha = 0,$$

And since the exponents $\frac{2}{\alpha} - 1, \frac{2}{\alpha} - 2 < 0$, this leads to a division by zero, or an undefined hessian matrix. However, as there are no other places in which $\frac{\partial^2 \Phi}{\partial x_i \partial x_j}$ is undefined, we have shown that the Hessian is well defined for all $\mathbf{x} \in \mathbb{R}^N - \mathbf{B}$.

Paul Lintilhac

2)

$$\frac{\partial \Phi}{\partial R_{t,i}} = 2 \left[\sum_{j=1}^N (R_{t,j})_+^\alpha \right]^{\frac{2}{\alpha}-1} (R_{t,i})_+^{\alpha-1}$$

We can re-write the dot product

$$\nabla \Phi(\mathbf{R}_{t-1}) * \mathbf{r}_t = \sum_{i=1}^N \frac{\partial \Phi}{\partial R_{t-1,i}} * \mathbf{r}_{t,i}$$

Using the fact that $r_{t,i} = L(\hat{y}_t, y_t) - L(y_{t,i}, y_t)$

$$= \sum_{i=1}^N 2 \left[\sum_{j=1}^N (R_{t-1,j})_+^\alpha \right]^{\frac{2}{\alpha}-1} (R_{t-1,i})_+^{\alpha-1} * (L(\hat{y}_t, y_t) - L(y_{t,i}, y_t))$$

Note that

$$\hat{y}_t = \frac{\sum_{i=1}^N w_{t,i} y_{t,i}}{\sum_{i=1}^N w_{t,i}}$$

And

$$w_{t,i} = (R_{t-1,i})_+^{\alpha-1}$$

Therefore we can write

$$\begin{aligned} & (L(\hat{y}_t, y_t) - L(y_{t,i}, y_t)) \\ &= L\left(\frac{\sum_{j=1}^N w_{t,j} y_{t,j}}{\sum_{j=1}^N w_{t,j}}, y_t\right) - L(y_{t,i}, y_t) \\ &\leq \frac{\sum_{j=1}^N w_{t,j} L(y_{t,j}, y_t)}{\sum_{j=1}^N w_{t,j}} - L(y_{t,i}, y_t) \\ & \text{(convexity of } L \text{ wrt first argument)} \\ &\leq L(y_{t,i}, y_t) - L(y_{t,i}, y_t) = 0 \\ & \text{(Cauchy - Schwartz)} \end{aligned}$$

Since the rest of the terms in the summand are all non-negative, this implies that

$$\nabla \Phi(\mathbf{R}_{t-1}) * \mathbf{r}_t \leq 0$$

Paul Lintilhac

3) We define matrices \mathbf{D} and \mathbf{P} as follows:

$$\{\mathbf{D}_{ii} : i \in \{1, \dots, N\}\} = (2\alpha - 2) \left[\sum_{k=1}^N (R_{t-1,k})_+^\alpha \right]^{\frac{2}{\alpha}-1} (R_{t-1,i})_+^{\alpha-2} \quad (\mathbf{D}_{ij} = 0 \text{ if } i \neq j)$$

$$\{\mathbf{P}_{ij} : i \in \{1, \dots, N\}\} = (4 - 2\alpha) \left[\sum_{k=1}^N (R_{t-1,k})_+^\alpha \right]^{\frac{2}{\alpha}-2} (R_{t-1,i})_+^{\alpha-1} (R_{t-1,j})_+^{\alpha-1}$$

Note that $\mathbf{r}_t^T \nabla^2 \phi(\mathbf{R}_{t-1}) \mathbf{r}_t = \mathbf{r}_t^T \mathbf{D} \mathbf{r}_t + \mathbf{r}_t^T \mathbf{P} \mathbf{r}_t$. Let us define $\mathbf{S} = \sum_{k=1}^N (R_{t-1,k})_+^\alpha$. Then

$$\begin{aligned} & \mathbf{r}_t^T \nabla^2 \phi(\mathbf{R}_{t-1}) \mathbf{r}_t \\ &= \left((2\alpha - 2) \mathbf{S}^{\frac{2}{\alpha}-1} \sum_i r_{t,i}^2 (R_{t-1,i})_+^{\alpha-2} + (4 - 2\alpha) \mathbf{S}^{\frac{2}{\alpha}-2} \sum_{i,j} r_{t,i} r_{t,j} (R_{t-1,i})_+^{\alpha-1} (R_{t-1,j})_+^{\alpha-1} \right) \end{aligned}$$

Note that \mathbf{P} is a positive semi-definite matrix, so the second term must be non-negative. However, notice that $\mathbf{S}^{\frac{2}{\alpha}-1} (\mathbf{R}_{t-1})_+^{\alpha-1} = \nabla \phi(\mathbf{R}_{t-1})$. From the previous question, we know that $\mathbf{S}^{\frac{2}{\alpha}-1} \sum_{i,j} r_{t,i} r_{t,j} (R_{t-1,i})_+^{\alpha-1} = \nabla \phi(\mathbf{R}_{t-1}) * \mathbf{r}_t \leq 0$. Thus we conclude that the second term must be 0:

$$\mathbf{r}_t^T \nabla^2 \phi(\mathbf{R}_{t-1}) \mathbf{r}_t = (2\alpha - 2) \mathbf{S}^{\frac{2}{\alpha}-1} \sum_i r_{t,i}^2 (R_{t-1,i})_+^{\alpha-2}$$

Note the additional term because elements along the diagonal need to be counted twice. Also note that

$$\begin{aligned} \mathbf{S}^{\frac{2}{\alpha}-1} \sum_i r_{t,i}^2 (R_{t-1,i})_+^{\alpha-2} &= \frac{\sum_i r_{t,i}^2 (R_{t-1,i})_+^{\alpha-2}}{\left(\sum_{k=1}^N (R_{t-1,k})_+^\alpha \right)^{\frac{\alpha-2}{\alpha}}} \\ &= \frac{\sum_i r_{t,i}^2 (R_{t-1,i})_+^{\alpha-2}}{\left(\sum_{k=1}^N (R_{t-1,k})_+^{\alpha-2} \right)^{\frac{\alpha}{\alpha-2}}} < r_{t,i}^2 (R_{t-1,k})_+^{\alpha-2} > \\ &= \frac{\sum_i r_{t,i}^2 (R_{t-1,i})_+^{\alpha-2}}{\left\| (R_{t-1,k})_+^{\alpha-2} \right\|_{\frac{\alpha}{\alpha-2}}} \end{aligned}$$

Now, we can apply Holder's inequality to conclude that

$$\mathbf{S}^{\frac{2}{\alpha}-1} \sum_i r_{t,i}^2 (R_{t-1,i})_+^{\alpha-2} \leq \left\| r_{t,i}^2 \right\|_{\frac{\alpha}{2}} = \left\| r_{t,i} \right\|_{\alpha}^2$$

Substituting these results back in, we get

$$\mathbf{r}_t^T \nabla^2 \phi(\mathbf{R}_{t-1}) \mathbf{r}_t \leq (2\alpha - 2) \left\| r_{t,i} \right\|_{\alpha}^2$$

Paul Lintilhac

4)

We can use Taylor's formula with a remainder up to the second order. Note that although we have shown that Φ is not necessarily differentiable more than 2 times, we did not show that it is necessarily non-differentiable (for example we could have $\alpha > 3$), and it could still be the case that $\Phi^{(3)}$ is defined.

$$\Phi(\mathbf{R}_t) = \Phi(\mathbf{R}_{t-1}) + \mathbf{r}_t^T \nabla \Phi(\mathbf{R}_{t-1}) + \frac{1}{2} \mathbf{r}_t^T \nabla^2 \Phi(\mathbf{R}_{t-1}) \mathbf{r}_t + \frac{1}{6} \Phi^{(3)}(\xi) \mathbf{r}_t^3$$

where $\xi \in [\mathbf{R}_{t-1}, \mathbf{R}_t]$

Note that any number in the range $[\mathbf{R}_{t-1}, \mathbf{R}_t]$ can be expressed as $\gamma \mathbf{R}_{t-1} + (1 - \gamma) \mathbf{R}_t$ where $\gamma \in [0, 1]$. Since by assumption $\gamma \mathbf{R}_{t-1} + (1 - \gamma) \mathbf{R}_t$ is not in \mathbf{B} , this implies that ξ is not in \mathbf{B} .

Note that this is important because the third derivative derivatives will include a multiplicative factor of

$$\left[\sum_{j=1}^N (\xi)_+^{\alpha} \right]^{\frac{2}{\alpha}-3}$$

and if $\xi \in \mathbf{B}$, this becomes undefined because of the negative exponent.

Therefore the second order approximation is:

$$\begin{aligned} \Rightarrow \Phi(\mathbf{R}_t) - \Phi(\mathbf{R}_{t-1}) &\approx \mathbf{r}_t^T \nabla \Phi(\mathbf{R}_{t-1}) + \frac{1}{2} \mathbf{r}_t^T \nabla^2 \Phi(\mathbf{R}_{t-1}) \mathbf{r}_t \\ &\leq 0 + (\alpha - 1) \|\mathbf{r}_t\|_{\alpha}^2 \end{aligned}$$

Where I have used the fact that $\mathbf{R}_t - \mathbf{R}_{t-1} = \mathbf{r}_t$.

5)

Suppose such a γ does exist. This means that there is some number $\xi \in [\mathbf{R}_{t-1}, \mathbf{R}_t]$ such that $\xi \in \mathbf{B}$. If we consider the fact that the total regret is an increasing function, this implies that $\mathbf{R}_{t-1} = \mathbf{0}$. If this is the case, then $\Phi(\mathbf{R}_{t-1}) = 0$, and plugging in to the result from the previous question, we get

$$\Phi(\mathbf{R}_t) \approx \mathbf{r}_t^T \nabla \Phi(\mathbf{R}_{t-1}) + \frac{1}{2} \mathbf{r}_t^T \nabla^2 \Phi(\mathbf{R}_{t-1}) \mathbf{r}_t \leq (\alpha - 1) \|\mathbf{r}_t\|_{\alpha}^2$$

Paul Lintilhac

6)

Now that we have a formula bounding $\Phi(\mathbf{R}_t) - \Phi(\mathbf{R}_{t-1})$. First, note that

$$\Phi(\mathbf{R}_0) = \left[\sum_{i=1}^N (R_{0,i})_+^\alpha \right]^{\frac{2}{\alpha}} = \left[\sum_{i=1}^N 0 \right]^{\frac{2}{\alpha}} = 0.$$

Then we can express $\Phi(\mathbf{R}_T)$ as a telescoping sum:

$$\begin{aligned} \Phi(\mathbf{R}_T) &= \Phi(\mathbf{R}_T) - \Phi(\mathbf{R}_0) = \sum_{t=0}^{T-1} \Phi(\mathbf{R}_{t+1}) - \Phi(\mathbf{R}_t) \\ &\leq \sum_{t=0}^{T-1} (\alpha - 1) \|\mathbf{r}_t\|_\alpha^2 = (\alpha - 1) \sum_{t=0}^{T-1} \|\mathbf{r}_t\|_\alpha^2 \end{aligned}$$

Now, in order to provide a (loose) upper bound, it would suffice to put an upper bound on $\|\mathbf{r}_t\|_\alpha^2$.

Note that

$$\|\mathbf{r}_t\|_\alpha^2 = \left(\sum_{i=1}^N \left(L(\hat{y}_t, y_t) - L(y_{t,i}, y_t) \right)^\alpha \right)^{\frac{2}{\alpha}}$$

Now, note that since $L(\hat{y}_t, y_t) \in [0, M]$, we know that $L(\hat{y}_t, y_t) - L(y_{t,i}, y_t) \leq M$.

Thus

$$\|\mathbf{r}_t\|_\alpha^2 \leq \left(\sum_{i=1}^N M^\alpha \right)^{\frac{2}{\alpha}} = (NM^\alpha)^{\frac{2}{\alpha}} = N^{\frac{2}{\alpha}} M^\alpha$$

Putting everything together, we have that

$$\Phi(\mathbf{R}_T) \leq (\alpha - 1)(T - 1) N^{\frac{2}{\alpha}} M^\alpha$$

Paul Lintilhac

7)

Showing the lower bound is straightforward. First, let me define $i^* = \underset{i \in \{1, \dots, N\}}{\operatorname{argmin}} \lambda(y_{T,i})$

Where $\lambda(y_{T,i}) = \sum_{t=1}^T L(y_{t,i}, y_t)$

Then

$$\begin{aligned} \Phi(\mathbf{R}_T) &= \left[\sum_{i=1}^N (R_{T,i})_+^\alpha \right]^{\frac{2}{\alpha}} = \left[\sum_{i=1}^N \left(\sum_t r_{t,i} \right)_+^\alpha \right]^{\frac{2}{\alpha}} = \left[\sum_{i=1}^N \left(\sum_t L(\widehat{y}_t, y_t) - L(y_{t,i}, y_t) \right)_+^\alpha \right]^{\frac{2}{\alpha}} \\ &= \left[\sum_{i=1}^N \left(\lambda(\widehat{y}_T) - \lambda(y_{T,i}) \right)_+^\alpha \right]^{\frac{2}{\alpha}} \\ &\geq \left[\sum_{i=1}^N \left(\lambda(\widehat{y}_T) - \lambda(y_{T,i^*}) \right)_+^\alpha \right]^{\frac{2}{\alpha}} \end{aligned}$$

Since $\lambda(y_{T,i^*}) \leq \lambda(y_{T,i})$ for all i

$$\geq \left[\left(\lambda(\widehat{y}_T) - \lambda(y_{T,i^*}) \right)_+^\alpha \right]^{\frac{2}{\alpha}}$$

Since $\left(\lambda(\widehat{y}_T) - \lambda(y_{T,i^*}) \right)_+^\alpha \geq 0$

$$\geq \left[\left(\lambda(\widehat{y}_T) - \lambda(y_{T,i^*}) \right)^\alpha \right]^{\frac{2}{\alpha}}$$

Since $\left(\lambda(\widehat{y}_T) - \lambda(y_{T,i^*}) \right)_+^\alpha \geq \left(\lambda(\widehat{y}_T) - \lambda(y_{T,i^*}) \right)^\alpha$

$$= \left(\lambda(\widehat{y}_T) - \lambda(y_{T,i^*}) \right)^2 = (R_T)^2$$

8) Given the above two bounds, we can make a straightforward combination of them:

$$(R_T)^2 \leq \Phi(\mathbf{R}_T) \leq (\alpha - 1)(T - 1)N^{\frac{2}{\alpha}} M^\alpha$$

$$\Rightarrow R_T \leq \sqrt{(\alpha - 1)(T - 1)N^{\frac{2}{\alpha}} M^\alpha}$$