

1.1 Bounded Difference (McDiarmid's) Inequality

At this point, we have the tools to control the size of expected supremum $\mathbb{E}\|P_n - P\|_{\mathcal{F}}$. Result in this section allows one to show that the random variable $\|P_n - P\|_{\mathcal{C}}$ concentrates tightly around its expectation.

Let $\mathbb{Z}(x_1, \dots, x_n)$ be a function of n variables. Assume that it satisfies the following assumption:

Assumption (Bounded difference). For all $1 \leq j \leq n$,

$$\sup_{x_1, \dots, x_j, x'_j, \dots, x_n} |\mathbb{Z}(x_1, \dots, x_j, \dots, x_n) - \mathbb{Z}(x_1, \dots, x'_j, \dots, x_n)| \leq c_j.$$

In essence, the bounded difference condition controls the fluctuations of \mathbb{Z} with respect to each variable separately.

Theorem 1. Let X_1, \dots, X_n be i.i.d. random variables. Assume that $\mathbb{Z}(X_1, \dots, X_n)$ satisfies the bounded difference assumption. Then

$$\Pr(|\mathbb{Z}(X_1, \dots, X_n) - \mathbb{E}\mathbb{Z}(X_1, \dots, X_n)| \geq t) \leq 2e^{-\frac{t^2}{2 \sum_{j=1}^n c_j^2}}.$$

Remark 1. A more careful argument (for instance, given in the original paper by McDiarmid) allows to obtain the tail bound $2e^{-\frac{2t^2}{\sum_{j=1}^n c_j^2}}$.

Example 1. Let $\mathbb{Z}(X_1, \dots, X_n) = \sum_{j=1}^n X_j$. Then

$$\mathbb{Z}(X_1, \dots, X_j, \dots, X_n) - \mathbb{Z}(X_1, \dots, X'_j, \dots, X_n) = X_j - X'_j.$$

Hence, if $|X_j| \leq M \implies |X_j - X'_j| \leq 2M$, and the bounded difference inequality becomes Hoeffding's inequality.

Example 2. This is one of the main examples we are interested in. Let

$$\mathbb{Z} = \|P_n - P\|_{\mathcal{C}} = \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{j=1}^n I\{X_j \in C\} - P(C) \right|.$$

Then \mathbb{Z} satisfies the bounded difference properly with $c_j = \frac{1}{n}$. Indeed,

$$\begin{aligned} & \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{j=1}^n I\{X_j \in C\} - P(C) + \frac{1}{n} I\{X'_j \in C\} - \frac{1}{n} I\{X'_j \in C\} \right| \\ & \leq \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{j \neq i} I\{X_j \in C\} + I\{X'_j \in C\} - P(C) \right| + \left| \frac{I\{X_j \in C\}}{n} - \frac{I\{X'_j \in C\}}{n} \right|. \end{aligned}$$

Proof of the Theorem. The proof is based on the *martingale decomposition technique*. To this end, we will find random variables V_1, \dots, V_n such that

- (1) $\mathbb{Z}(X_1, \dots, X_n) - \mathbb{E}\mathbb{Z}(X_1, \dots, X_n) = \sum_{j=1}^n V_j$;
- (2) $\mathbb{E}[V_j | V_1, \dots, V_{j-1}] = 0$;
- (3) $V_j = V_j(X_1, \dots, X_j)$ (in other words, V_j is a function of X_1, \dots, X_j only).

Recall the definition of a martingale sequence:

Definition 1 (Martingale). A sequence $(Y_j, \mathcal{F}_j)_{j \geq 1}$, where X_j is a random variable and \mathcal{F}_j is a σ -algebra, is a martingale if

- (1) $\mathbb{E}|Y_j| < \infty$;
- (2) Y_j is \mathcal{F}_j – measurable;
- (3) $\mathbb{E}(Y_j | \mathcal{F}_{j-1}) = Y_{j-1}$.

Note that the partial sums $\sum_{j=1}^k V_j$ form a martingale, assuming that V_j 's satisfy the requirements stated above.

We will define V_1, \dots, V_n as follows: let (Y_1, \dots, Y_n) be an independent copy of (X_1, \dots, X_n) , and set

$$V_j = \mathbb{E}_Y[\mathbb{Z}(X_1, \dots, X_j, Y_{j+1}, \dots, Y_n)] - \mathbb{E}_Y[\mathbb{Z}(X_1, \dots, X_{j-1}, Y_j, \dots, Y_n)].$$

From this definition, it is clear that $|V_j| \leq c_j$ by the bounded difference property of \mathbb{Z} . Moreover, note that V_j can be equivalently written as

$$V_j = \mathbb{E}[\mathbb{Z}(X_1, \dots, X_n) | X_1, \dots, X_j] - \mathbb{E}[\mathbb{Z}(X_1, \dots, X_n) | X_1, \dots, X_{j-1}].$$

It is easy to check that $\sum_{j=1}^n V_j = \mathbb{Z} - \mathbb{E}\mathbb{Z}$ and that V_j depends only on X_1, \dots, X_j . To show the second property, note that by the “tower property” of conditional expectations,

$$\begin{aligned} \mathbb{E}[\mathbb{E}\mathbb{Z} | X_1, \dots, X_j] - \mathbb{E}[\mathbb{Z}(X_1, \dots, X_{j-1}) | X_1, \dots, X_{j-1}] \\ = \mathbb{E}[\mathbb{Z} | X_1, \dots, X_{j-1}] - \mathbb{E}[\mathbb{Z} | X_1, \dots, X_{j-1}] = 0. \end{aligned}$$

The rest of the proof proceeds in a way that is similar to the proofs of Hoeffding's and Bernstein's inequalities (that is, estimating the moment generating function). Let $\lambda > 0$ (to be chosen later), and note that

$$\begin{aligned} \Pr(\mathbb{Z} - \mathbb{E}\mathbb{Z} > t) &= \Pr(\lambda(\mathbb{Z} - \mathbb{E}\mathbb{Z}) \geq \lambda t) \\ &= \Pr(e^{\lambda(\mathbb{Z} - \mathbb{E}\mathbb{Z})} \geq e^{\lambda t}) \leq e^{-\lambda t} \mathbb{E}e^{\lambda(\mathbb{Z} - \mathbb{E}\mathbb{Z})} \\ &= e^{-\lambda t} \mathbb{E}e^{\lambda \sum_{j=1}^n V_j}. \end{aligned}$$

Next,

$$\begin{aligned}
\mathbb{E}e^{\lambda \sum_{j=1}^n V_j} &= \mathbb{E}\mathbb{E}\left[e^{\lambda \sum_{j=1}^n V_j} | X_1, \dots, X_{n-1}\right] = \mathbb{E}\mathbb{E}\left[e^{\lambda \sum_{j=1}^{n-1} V_j} e^{\lambda V_n} | X_1, \dots, X_{n-1}\right] \\
&= \mathbb{E}\left(e^{\lambda \sum_{j=1}^{n-1} V_j} \mathbb{E}\left[e^{\lambda V_n} | X_1, \dots, X_{n-1}\right]\right) \leq e^{\frac{\lambda c_n^2}{2}} \mathbb{E}e^{-\lambda \sum_{j=1}^{n-1} V_j} \\
&= e^{\frac{\lambda c_n^2}{2}} \mathbb{E}\mathbb{E}\left[e^{\lambda \sum_{j=1}^{n-2} V_j} e^{\lambda V_{n-1}} | X_1, \dots, X_{n-2}\right] \\
&= e^{\frac{\lambda c_n^2}{2}} \mathbb{E}\left(e^{\lambda \sum_{j=1}^{n-2} V_j} \mathbb{E}\left[e^{\lambda V_{n-1}} | X_1, \dots, X_{n-2}\right]\right) \\
&\leq e^{\frac{\lambda^2 c_n^2}{2}} e^{\frac{\lambda^2 c_{n-1}^2}{2}} \mathbb{E}e^{\lambda \sum_{j=1}^{n-2} V_j} \leq \dots \leq e^{\frac{\lambda^2}{2} \sum_{j=1}^n c_j^2}.
\end{aligned}$$

We used the property that $|V_j| \leq c_j$ almost surely, and the fact that bounded zero-mean random variables are sub-Gaussian, implying that $\mathbb{E}\left[e^{\lambda V_j} | X_1, \dots, X_{j-1}\right] \leq e^{\lambda^2 c_j^2 / 2}$. Hence,

$$\Pr(\mathbb{Z} - \mathbb{E}\mathbb{Z} \geq t) \leq e^{-\lambda t + \frac{\lambda^2}{2} \sum_{j=1}^n c_j^2}$$

The minimum of $f(\lambda) = -\lambda t + \frac{\lambda^2}{2} \sum_{j=1}^n c_j^2$ is attained for $\lambda_* = \frac{t}{\sum_{j=1}^n c_j^2}$, and the final result follows.

The inequality for the lower tail $\Pr(\mathbb{Z} - \mathbb{E}\mathbb{Z} \leq -t)$ is obtained from the same argument applied to $-Z$ instead of Z . \square

1.2 Generalization error bounds for Adaboost

We are almost ready to go back to the first weeks of the course and establish a bound on the generalization error of AdaBoost. To this end, we will need one more technical result called “Talagrand’s contraction inequality” that is stated below. Its (very nice!) proof is given in the end of these notes.

We will say that $\varphi : \mathbb{R} \mapsto \mathbb{R}$ is a *contraction* whenever $|\varphi(x) - \varphi(y)| \leq |x - y|$ for all $x, y \in \mathbb{R}$.

Theorem 2 (M. Talagrand). Assume that Φ be a convex function that is nondecreasing on $[0, \infty)$, and let φ_j be contractions for $j = 1, \dots, n$ such that $\varphi_j(0) = 0$. Then for any $T \subseteq \mathbb{R}^n$

$$\mathbb{E}\Phi\left(\frac{1}{2} \sup_{t \in T} \left| \sum_{j=1}^n \varepsilon_j \varphi(t_j) \right| \right) \leq \mathbb{E}\Phi\left(\sup_{t \in T} \left| \sum_{j=1}^n \varepsilon_j t_j \right| \right).$$

Here, $t = (t_1, \dots, t_n)$. Before we proceed to the proof, let’s look at a useful example.

Example 3. Let $\mathcal{F} = \{f : S \mapsto \mathbb{R}\}$, and let φ be a contraction. Define $\varphi \circ \mathcal{F} := \{\varphi(f), f \in \mathcal{F}\}$. Finally, given $x_1, \dots, x_n \in S$, let $T := \{f(x_1), \dots, f(x_n), f \in \mathcal{F}\}$. Then

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j \varphi(f(x_j)) \right| \leq 2 \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j f(x_j) \right|.$$

For instance, of $\varphi(x) = x^2$ and $|f| \leq U$ for $f \in \mathcal{F}$, then $\varphi(x) = \frac{\varphi(x)}{2U}$ is a contraction on $[-U, U]$. Hence,

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j f^2(x_j) \right| \leq 4U \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j f(x_j) \right|$$

In particular, we can replace x_1, \dots, x_n by random variables X_1, \dots, X_n , and get that

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j \phi(f(X_j)) \right| \leq 2 \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j f(X_j) \right|$$

Now we proceed with our main goal – controlling the generalization error of AdaBoost. Let \mathcal{F} be the “base class” – class of binary classifiers of VC dimension V , and let

$$\hat{g}_T = \frac{\sum_{j=1}^T \alpha_j f_j}{\sum_{j=1}^T \alpha_j}$$

be the output of Adaboost after T steps; note that $\hat{g}_T \in \text{co}(\mathcal{F})$, the convex hull of \mathcal{F} .

Theorem 3. With probability $\geq 1 - e^{-t}$, for any $\theta > 0$

$$\Pr(Y \neq \text{sign}(\hat{g}_T(X))) \leq \frac{1}{n} \sum_{j=1}^n \mathbb{I}\{Y_j \hat{g}_T(X_j) \leq \theta\} + K \left(\frac{1}{\theta} \sqrt{\frac{V}{n}} + \sqrt{\frac{t}{n}} \right)$$

where K is an absolute constant.

Proof. Let

$$\varphi_\theta(x) = \begin{cases} 1, & x \leq 0 \\ 1 - \frac{x}{\theta}, & 0 \leq x \leq \theta, \\ 0, & x \geq \theta \end{cases}$$

and note that φ_θ is Lipschitz continuous with Lipschitz constant $\frac{1}{\theta}$. Moreover, it is easy to see that $\mathbb{I}\{x \leq 0\} \leq \varphi_\theta(x) \leq \mathbb{I}\{x \leq \theta\}$ for all $x \in \mathbb{R}$. Consequently,

$$\begin{aligned} \Pr(Y \hat{g}_T(X) \leq 0) &= \mathbb{E} \mathbb{I}\{Y \hat{g}_T(X) \leq 0\} \leq \mathbb{E} \varphi_\theta(Y \hat{g}_T(X)) \\ &= \mathbb{E} \varphi_\theta(Y \hat{g}_T(X)) \pm \left(\frac{1}{n} \sum_{j=1}^n \varphi_\theta(Y_j \hat{g}_T(X_j)) \right) \\ &\leq \frac{1}{n} \sum_{j=1}^n \mathbb{I}\{Y_j \hat{g}_T(X_j) \leq \theta\} + \underbrace{\sup_{g \in \text{co}(\mathcal{F})} \left| \frac{1}{n} \sum_{j=1}^n \varphi_\theta(Y_j g(X_j)) - \mathbb{E} \varphi_\theta(Y g(X)) \right|}_{Z((X_1, Y_1), \dots, (X_n, Y_n))}. \end{aligned}$$

Notice that Z satisfies the bounded difference property with $c_j = \frac{1}{n}$. Hence, the bounded difference inequality implies that

$$Z \leq \mathbb{E} Z + \sqrt{\frac{t}{n}} \text{ with probability } \geq 1 - e^{-2t}.$$

It remains to estimate $\mathbb{E}Z$. To this end, we apply the symmetrization and contraction inequalities to get that

$$\begin{aligned} \mathbb{E} \sup_{g \in \text{co}(\mathcal{F})} \left| \frac{1}{n} \sum_{j=1}^n (\varphi_n(Y_j g(X_j))) - \mathbb{E} \varphi_\theta(Y_j g(X_j)) \right| &\leq \frac{2}{\theta} \mathbb{E} \sup_{g \in \text{co}(\mathcal{F})} \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j \varphi_\theta(Y_j g(X_j)) \right| \\ &\leq \frac{4}{\theta} \mathbb{E} \sup_{g \in \text{co}(\mathcal{F})} \left| \frac{1}{n} \sum_{j=1}^n \varepsilon_j Y_j g(X_j) \right|. \end{aligned}$$

Since the maximum of a linear function over a convex set is always attained at one of the extreme points,

$$\mathbb{E} \sup_{g \in \text{co}(\mathcal{F})} \left| \frac{1}{n} \sum_{j=1}^n \tilde{\varepsilon}_j g(X_j) \right| = \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \tilde{\varepsilon}_j f(X_j) \right|,$$

where $\tilde{\varepsilon}_j = \varepsilon_j Y_j$, $j = 1, \dots, n$ is a sequence of Rademacher random variables since ε_j 's and Y_j 's are independent. It implies that

$$\begin{aligned} \mathbb{E} \sup_{g \in \text{co}(\mathcal{F})} \left| \frac{1}{n} \sum_{j=1}^n (\varphi_n(Y_j g(X_j))) - \mathbb{E} \varphi_\theta(Y_j g(X_j)) \right| &\leq \frac{4}{\theta} \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \tilde{\varepsilon}_j f(X_j) \right| \\ &\leq \frac{K}{\theta} \sqrt{\frac{V(\mathcal{F})}{n}}, \end{aligned}$$

where K is an absolute constant; the last inequality follows from Dudley's entropy integral bound and Haussler's theorem. \square

Finally, we will address the size of the term $\frac{1}{n} \sum_{j=1}^n \mathbb{I}\{Y_j \hat{g}_T(X_j) \leq \theta\}$. Assume that γ -weak learnability assumption holds, meaning that for some $\gamma > 0$

$$\forall \omega_1, \dots, \omega_n \geq 0, \sum_{j=1}^n \omega_j = 1, \exists f \in \mathcal{F}, \text{ such that } \sum_{j=1}^n \omega_j \mathbb{I}\{Y_j \neq f(X_j)\} \leq 1/2 - \gamma.$$

Theorem 4. Let \hat{g}_T be the output of AdaBoost after T steps, and assume that γ -weak learnability assumption holds (for some, possibly unknown, $\gamma > 0$). Then $\forall \theta \in (0, 1)$,

$$\frac{1}{n} \sum_{j=1}^n \mathbb{I}\{Y_j \hat{g}_T(X_j) \leq \theta\} \leq 2^T \left(\sqrt{(1 - 2\gamma)^{1-\theta} (1 + 2\gamma)^{1+\theta}} \right)^T.$$

Remark 2. As $\theta \rightarrow 0$, the right-hand side of the inequality converges to $\left(\sqrt{(\frac{1}{2} - \gamma)(\frac{1}{2} + \gamma)} \right)^T$ that goes to 0 exponentially fast with T . Hence, for θ small enough, it is bounded by $C(\gamma, \theta)^T$ for some constant $C(\gamma, \theta) < 1$.

Exercise. Prove that $C(\gamma, \theta) < 1$ if $\theta < \gamma$.

Sketch of the proof. Assume that $Y\hat{g}_T(X) \leq \theta$, then

$$\exp\left(-Y \sum_{j=1}^T \alpha_j f_j(X) + \theta \sum_{j=1}^T \alpha_j\right) \geq 1.$$

Next,

$$\begin{aligned} \frac{1}{n} \sum_{j=1}^n \mathbb{I}\{Y_j \hat{g}_T(X_j) \leq \theta\} &= \frac{1}{n} \sum_{j=1}^n \mathbb{I}\{Y_j \hat{g}_T(X_j) - \theta \leq 0\} \\ &\leq \frac{1}{n} \sum_{j=1}^n \exp\left(\theta \sum_{i=1}^T \alpha_i - Y_j \sum_{i=1}^T \alpha_i f_i(X_j)\right) \\ &= \frac{1}{n} \exp\left(\theta \sum_{i=1}^T \alpha_i\right) \sum_{j=1}^n \exp\left(-Y_j \sum_{i=1}^T \alpha_i f_i(X_j)\right). \end{aligned}$$

The rest of the proof (estimation of the term $\sum_{j=1}^n \exp\left(-Y_j \sum_{i=1}^T \alpha_i f_i(X_j)\right)$) proceeds as before, see Theorem 1 in the notes for Week 4. \square

1.3 Proof of Talagrand's Contraction Inequality

1. First, we will show that for an *arbitrary map* $A : T \subset \mathbb{R}^n \mapsto \mathbb{R}$,

$$\mathbb{E}\Phi\left(\sup_{t \in T} \left[A(t) + \sum_{j=1}^n \varepsilon_j \varphi(t_j)\right]\right) \leq \mathbb{E}\Phi\left(\sup_{t \in T} \left[A(t) + \sum_{j=1}^n \varepsilon_j t_j\right]\right)$$

2. It is sufficient to consider the case $n = 1$. Indeed, if

$$\mathbb{E}\Phi\left(\sup_{t \in T} [A(t) + \varepsilon \varphi(t_n)]\right) \leq \mathbb{E}\Phi\left(\sup_{t \in T} [A(t) + \varepsilon t_n]\right)$$

then

$$\begin{aligned} &\mathbb{E}_{\varepsilon_1, \dots, \varepsilon_{n-1}} \mathbb{E}_{\varepsilon_n} \Phi\left(\sup_{t \in T} \left[\underbrace{A(t) + \sum_{j=1}^{n-1} \varepsilon_j \varphi(t_j)}_{\hat{A}(t)} + \varepsilon_n \varphi(t_n)\right]\right) \\ &\leq \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_{n-1}} \mathbb{E}_{\varepsilon_n} \Phi\left(\sup_{t \in T} \left[A(t) + \sum_{j=1}^{n-1} \varepsilon_j \varphi_j(t) + \varepsilon_n t_n\right]\right) \\ &\leq \dots \leq \mathbb{E}\Phi\left(\sup_{t \in T} \left[A(t) + \sum_{j=1}^n \varepsilon_j t_j\right]\right). \end{aligned}$$

Recall that we want to show that $\mathbb{E}\Phi(\sup_{t \in T} [A(t) + \varepsilon\varphi(t_n)]) \leq \mathbb{E}\Phi(\sup_{t \in T} [A(t) + \varepsilon t_n])$. This is equivalent to

$$\begin{aligned} \frac{1}{2}\Phi\left(\sup_{t \in T} [A(t) - \varphi(t_n)]\right) + \frac{1}{2}\Phi\left(\sup_{t \in T} [A(t) + \varphi(t_n)]\right) \\ \leq \frac{1}{2}\Phi\left(\sup_{t \in T} [A(t) - t_n]\right) + \frac{1}{2}\Phi\left(\sup_{t \in T} [A(t) + t_n]\right) \end{aligned}$$

Let $s_1 := A(t)$, $s_2 := t_n$, and $\tilde{T} = \{(s_1, s_2), t \in T\}$. We then need to show that

$$\begin{aligned} \Phi\left(\sup_{s_1, s_2 \in \tilde{T}} [s_1 - \varphi(s_2)]\right) + \Phi\left(\sup_{s_1, s_2 \in \tilde{T}} [s_1 + \varphi(s_2)]\right) \\ \leq \Phi\left(\sup_{s_1, s_2 \in \tilde{T}} (s_1 - s_2)\right) + \Phi\left(\sup_{s_1, s_2 \in \tilde{T}} (s_1 + s_2)\right) \end{aligned}$$

Let \bar{t}_1, \bar{t}_2 be such that $\sup_{s_1, s_2} (t_1 + \varphi(t_2)) = \bar{t}_1 + \varphi(\bar{t}_2)$, and \bar{s}_1, \bar{s}_2 — such that $\sup_{s_1, s_2} (s_1 - \varphi(s_2)) = \bar{s}_1 - \varphi(\bar{s}_2)$. We will consider several cases:

(a) Assume that $\bar{t}_2 \geq 0, \bar{s}_2 \geq 0, \bar{t}_2 \geq \bar{s}_2$. We will show that

$$\Phi(\bar{t}_1 + \varphi(\bar{t}_2)) + \Phi(\bar{s}_1 - \varphi(\bar{s}_2)) \leq \Phi(\bar{t}_1 + \bar{t}_2) + \Phi(\bar{s}_1 - \bar{s}_2),$$

which suffices as the right-hand side is clearly at most $\Phi(\sup_{s_1, s_2 \in \tilde{T}} (s_1 - s_2)) + \Phi(\sup_{s_1, s_2 \in \tilde{T}} (s_1 + s_2))$. Let $a = \bar{t}_1 + \varphi(\bar{t}_2)$, $b = \bar{t}_1 + \bar{t}_2$, $c = \bar{s}_1 - \bar{s}_2$ and $d = \bar{s}_1 - \varphi(\bar{s}_2)$, and note that $a \leq b, c \leq d$ since φ is a contraction. Next,

$$b - a = \bar{t}_2 - \varphi(\bar{t}_2) \geq \bar{s}_2 - \varphi(\bar{s}_2) = d - c,$$

hence $\varphi(\bar{t}_2) - \varphi(\bar{s}_2) \leq \bar{t}_2 - \bar{s}_2$. Moreover,

$$a = \bar{t}_1 + \varphi(\bar{t}_2) \geq \bar{s}_1 + \varphi(\bar{s}_2) \geq \bar{s}_1 - \bar{s}_2 = c.$$

We used the definition of \bar{t}_1 and \bar{t}_2 above, together with the fact that ϕ is a contraction. As Φ is nondecreasing and convex, these inequalities imply that the increment of Φ over $[a, b]$ is larger than the increment over $[c, d]$, or

$$\Phi(b) - \Phi(a) \geq \Phi(d) - \Phi(c).$$

The case $\bar{t}_2 \geq 0, \bar{s}_2 \geq 0, \bar{t}_2 \leq \bar{s}_2$ is similar after renaming the variables $(\bar{t}_1, \bar{t}_2) \leftrightarrow (\bar{s}_1, \bar{s}_2)$, and replacing φ by $-\varphi$.

(b) The case $\bar{t}_2 \leq 0, \bar{s}_2 \leq 0$ can be reduced to the previous one by using the transformation $(\bar{t}_1, \bar{t}_2) \mapsto (\bar{t}_1, -\bar{t}_2)$, $(\bar{s}_1, \bar{s}_2) \mapsto (\bar{s}_1, -\bar{s}_2)$ and $\tilde{\varphi}(x) = \varphi(-x)$.

(c) Finally, in the case $\bar{t}_2 \geq 0, \bar{s}_2 \leq 0$, observe that

$$\varphi(t_2) \leq t_2, \quad -\varphi(s_2) \leq -s_2,$$

hence monotonicity of Φ implies that

$$\Phi(\bar{t}_1 + \varphi(\bar{t}_2)) + \Phi(\bar{s}_1 - \varphi(\bar{s}_2)) \leq \Phi(\bar{t}_1 + \bar{t}_2) + \Phi(\bar{s}_1 - \bar{s}_2).$$

3. Let $a_+ := \max(a, 0)$. To finish the proof, note that

$$\begin{aligned} \mathbb{E}\Phi\left(\frac{1}{2}\sup_{t\in T}\left|\sum_{j=1}^n\varepsilon_j\varphi(t_j)\right|\right) &= \mathbb{E}\Phi\left(\frac{1}{2}\left(\sup_{t\in T}\sum_{j=1}^n\varepsilon_j\varphi(t_j)\right)_+ + \frac{1}{2}\left(\sup_{t\in T}\sum_{j=1}^n(-\varepsilon_j)\varphi(t_j)\right)_+\right) \\ &\leq \frac{1}{2}\mathbb{E}\Phi\left(\left(\sup_{t\in T}\sum_{j=1}^n\varepsilon_j\varphi(t_j)\right)_+\right) + \frac{1}{2}\mathbb{E}\Phi\left(\left(\sup_{t\in T}\sum_{j=1}^n(-\varepsilon_j)\varphi(t_j)\right)_+\right) \\ &= \mathbb{E}\Phi\left(\left(\sup_{t\in T}\sum_{j=1}^n\varepsilon_j\varphi(t_j)\right)_+\right) \end{aligned}$$

As the function $\tilde{\Phi}(a) := \Phi(a_+)$ is still non-decreasing and convex,

$$\mathbb{E}\Phi\left(\left(\sup_{t\in T}\sum_{j=1}^n\varepsilon_j\varphi(t_j)\right)_+\right) \leq \mathbb{E}\Phi\left(\sup_{t\in T}\left(\sum_{j=1}^n\varepsilon_j t_j\right)_+\right) \leq \mathbb{E}\Phi\left(\sup_{t\in T}\left|\sum_{j=1}^n\varepsilon_j t_j\right|\right),$$

where the last inequality follows from the monotonicity of Φ .