1) We have, by definition, for $i = 1, ..., n$.

$$\mathbb{1}_{\{X_i \leq t\}} = \begin{array}{ll} 1 \text{ if } X_i \leq t & \text{proba } \mathbb{P}(X_i \leq t) = F(t) \\ 0 \text{ if } X_i > t & \text{proba } 1 - F(t) \end{array}$$

So they are Bernoulli r.v. with parameter $F(t)$.
They are independent because the $X_i$'s are independent.

2) Apply Hoeffding's inequality, we get, for every $\varepsilon > 0$

$$\mathbb{P}\left( \left| \underbrace{\frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\{X_i \leq t\}}}_{= \hat{F}_n(t)} - F(t) \right| \geq \varepsilon \right) \leq 2 \exp\left(-2\varepsilon^2 n\right)$$

Since the right-hand side does not depend on $t$, we may write

$$\sup_{t \in \mathbb{R}} \mathbb{P}\left( |\hat{F}_n(t) - F(t)| \geq \varepsilon \right) \leq \sup_{t \in \mathbb{R}} 2 \exp\left(-2\varepsilon^2 n\right)$$
$$= 2 \exp\left(-2\varepsilon^2 n\right).$$

3) We always have

$$|\hat{F}_n(t) - F(t)| \leq \sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)|,$$

So

$$\mathbb{P}\left( |\hat{F}_n(t) - F(t)| \geq \varepsilon \right) \leq \mathbb{P}\left( \sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)| \geq \varepsilon \right).$$

We can take the supremum on $t$ in the left-hand side, the right-hand side does not depend on $t$! (I know, it's surprising, but observe that for any quantity $A$, $\sup_t A(t)$ can also be written $\sup_s A(s)$ or with whatever letter you want).

So $$\sup_{t \in \mathbb{R}} \mathbb{P}\left( |\hat{F}_n(t) - F(t)| \geq \varepsilon \right) \leq \mathbb{P}\left( \sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)| \geq \varepsilon \right),$$
and thus inequality (3) implies inequality (2).

**4)** No and no. Proof:

- We can compute

$$\mathbb{E}\left[\overline{P_n}(t)\right] = \frac{1}{\varepsilon}\left(\mathbb{E}\left[\hat{f}_n(t+\varepsilon)\right] - \mathbb{E}\left[\hat{f}_n(t)\right]\right)$$

$$= \frac{1}{\varepsilon}\left(F(t+\varepsilon) - F(t)\right) \quad \text{because } \hat{f}_n(e) \text{ is an}$$
$$\text{unbiased estimator}$$

this is "close" to $f(t)$ for $\varepsilon$ small, but
of $F(a)$ for all $a$.

not equal to it (in general).

So $\lim\limits_{n \to \infty} \mathbb{E}\left[\overline{P_n}(t)\right] \neq F'(t) = p(t)$. $\overline{P_n}(t)$ is Asymptotically **biased**.

- We know $\hat{f}_n(a)$ is a consistent estimator of $F(a)$, so

$$\hat{f}_n(t+\varepsilon) \xrightarrow[n \to \infty]{\mathbb{P}} F(t+\varepsilon), \quad \hat{f}_n(t) \xrightarrow[n \to \infty]{\mathbb{P}} F(t),$$

and thus $\overline{P_n}(t) \xrightarrow[n \to \infty]{\mathbb{P}} \frac{1}{\varepsilon}\left(F(t+\varepsilon) - F(t)\right)$

For the same reason ——— this is not $= F'(t)$.

$$\boxed{\text{So } \overline{P_n}(t) \text{ is not consistent}}$$

**5)** We can still write, for any $n, t$

$$\mathbb{E}\left[\hat{f}_n(t+\varepsilon_n)\right] = F(t+\varepsilon_n) \quad \text{because } \hat{f}_n \text{ is unbiased}.$$

we also have $\mathbb{E}\left[\hat{f}_n(t)\right] = F(t)$.

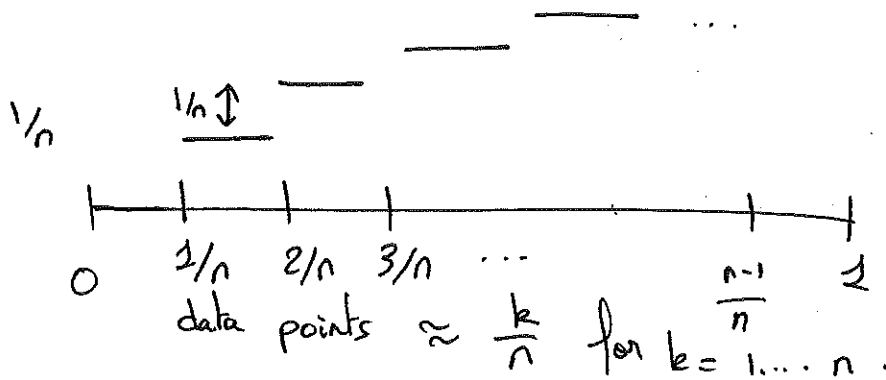So for any fixed $t$, for any $n$, we have

$$\mathbb{E}\left[\hat{P}_n(t)\right] = \frac{1}{\varepsilon_n}\left(F(t+\varepsilon_n) - F(t)\right).$$

Then, sending $n \to \infty$, since $\varepsilon_n \xrightarrow[n \to \infty]{} 0$, we have

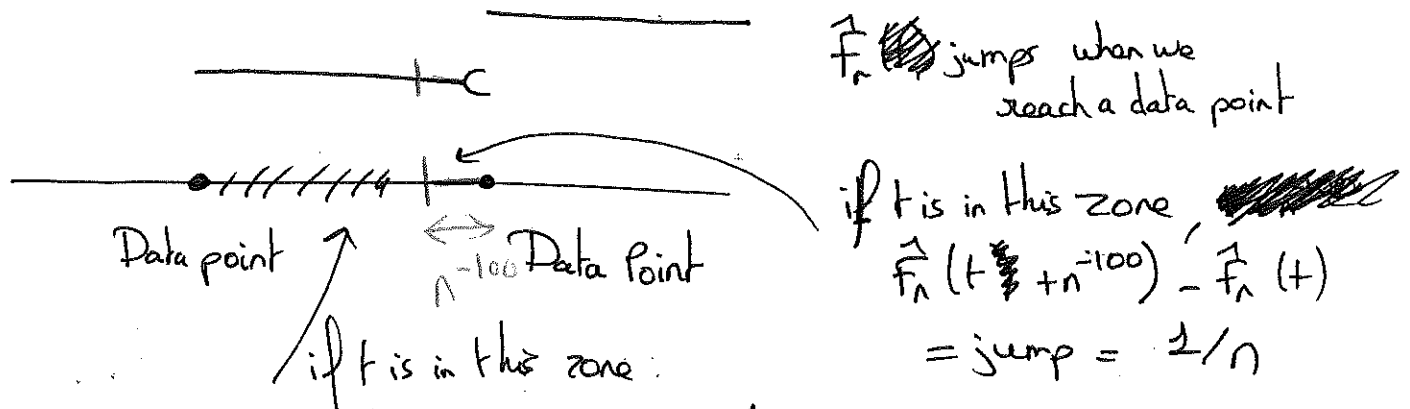$$\lim\limits_{n \to \infty} \mathbb{E}\left[\hat{P}_n(t)\right] = \lim\limits_{n \to \infty}\left(\frac{1}{\varepsilon_n}\left(F(t+\varepsilon_n) - F(t)\right)\right) = F'(t) = p(t)$$

So $\underline{\hat{P}_n(t) \text{ is asymptotically unbiased}}$.

**6]** Suppose we are sampling the uniform distribution on $[0,1]$.
Typically we will get $n$ data points that are (approximately) uniformly spread on $[0,1]$, so $\hat{F}_n(t)$ looks like

_____ 1
 _____
 _____ ...

$\frac{1}{n}$    $\frac{1}{n}\updownarrow$ _____

0   $\frac{1}{n}$ $\frac{2}{n}$ $\frac{3}{n}$ $\cdots$   $\frac{n-1}{n}$  1

data points $\approx \frac{k}{n}$ for $k = 1, \ldots, n$.

For $t$ arbitrary in $[0,1]$, $\hat{F}_n(t + n^{-100})$ and $\hat{F}_n(t)$ are equal, except if $t$ is near a data point (more precisely, except if $t$ is at distance $\leq n^{-100}$ on the left of a data point).

_____
_____ $+$ C

$\bullet$///////$+$   $+$      $\hat{F}_n$ jumps when we reach a data point

Data point ↑   $n^{-100}$ Data Point    if $t$ is in this zone,
                                         $\hat{F}_n(t + n^{-100}) - \hat{F}_n(t)$
if $t$ is in this zone                    $=$ jump $= \frac{1}{n}$
$$\hat{F}_n(t + n^{-100}) - \hat{F}_n(t) = 0.$$

So "most of the time", $\dfrac{\hat{F}_n(t + n^{-100}) - \hat{F}_n(t)}{n^{-100}} = 0.$

At some "exceptional" times, $\dfrac{\hat{F}_n(t + n^{-100}) - \hat{F}_n(t)}{n^{-100}} = \dfrac{\frac{1}{n}}{n^{-100}} = \dfrac{n^{99}}{\phantom{x}}$ Very large

In expectation, things work, because

Often $\times$ 0 $+$ Rarely $\times$ Very large $=$ OK $\rightarrow$ asymptotically unbiased, see Q.5.

However, "in probability", we only see the "often" part, and choosing $\varepsilon_n$ too small leads to an estimator $\hat{P}_n \xrightarrow[n \to \infty]{\mathbb{P}} 0$. Not consistent!

**7]** let t be fixed

We want to prove that $\hat{p}_n(t)$ is a consistent estimator of $p(t)$,

which means $\hat{p}_n(t) \xrightarrow[n\to\infty]{\mathbb{P}} p(t)$.

So we fix $\delta > 0$, and we want to prove

$$\mathbb{P}\left( \left| \hat{p}_n(t) - p(t) \right| > \delta \right) \xrightarrow[n\to\infty]{} 0 .$$

Let us write the definition of $\hat{p}_n$:

$$\hat{p}_n(t) = \frac{\hat{F}_n(t + n^{-1/4}) - \hat{F}_n(t)}{n^{-1/4}} .$$

We also have, since $p(t) = F'(t)$:

$$p(t) = \frac{F(t + n^{-1/4}) - F(t)}{n^{-1/4}} + h(n), \quad \text{with } h(n) \xrightarrow[n\to\infty]{} 0 ,$$

by definition of a derivative.

So

$$\left| \hat{p}_n(t) - p(t) \right| \leq \left( \frac{\left| \hat{F}_n(t + n^{-1/4}) - F(t + n^{-1/4}) \right|}{n^{-1/4}} \right.$$

$$\left. + \frac{\left| \hat{F}_n(t) - F(t) \right|}{n^{-1/4}} + \left| h(n) \right| \right) \qquad \text{by triangular inequality}$$

By inequality (1), we have, choosing $\varepsilon = \frac{\delta}{10} n^{-1/4}$:

$$\mathbb{P}\left( \left| \hat{F}_n(t) - F(t) \right| \geq \frac{\delta}{10} n^{-1/4} \right) \qquad \xrightarrow[\text{as } n\to\infty]{} 0 .$$

$$\leq 2 \exp\left( -2 \frac{\delta^2}{100} n^{-1/2} \cdot n \right) = 2\exp\left( -\frac{\delta^2}{50} n^{1/2} \right)$$

similarly, we have

$$\mathbb{P}\left( \left| \hat{F}_n(t + n^{-1/4}) - F_n(t + n^{-1/4}) \right| \geq \frac{\delta}{10} n^{-1/4} \right) \leq 2\exp\left( -\frac{\delta^2}{50} n^{1/2} \right).$$

So $\mathbb{P}\left( \frac{\left| \hat{F}_n(t + n^{-1/4}) - F(t + n^{-1/4}) \right|}{n^{-1/4}} + \frac{\left| \hat{F}_n(t) - F(t) \right|}{n^{-1/4}} \geq \frac{\delta}{5} \right)$

$$\underset{\text{bound}}{\overset{\text{Union}}{\longrightarrow}} \leq \textcircled{2} \, 2 \cdot \exp\left( -\frac{\delta^2}{50} n^{1/2} \right). \qquad \left( \frac{\delta}{5} = \frac{\delta}{10} + \frac{\delta}{10} \right)$$

Since $h(n) \longrightarrow 0$; it is less than $\frac{\delta}{5}$ for $n$ large enough.

In conclusion,

$$\frac{\left|\hat{F}_n\left(t+n^{-1/4}\right) - F\left(t+n^{-1/4}\right)\right|}{n^{-1/4}} + \frac{\left|\hat{F}_n(t) - F(t)\right|}{n^{-1/4}} + \left|h(n)\right|$$

is less than $\quad \delta/5 + \frac{\delta}{5} = \frac{2\delta}{5} < \delta,\quad$ for $n$ large enough,

with probability $\geq 1 - 4\exp\left(\frac{-\delta^2}{50} n^{1/2}\right), \longrightarrow \underset{as\ n\to\infty}{0}$ .

and thus $\mathbb{P}\left(\left|\hat{f}_n(t) - p(t)\right| \geq \delta\right) \xrightarrow[n\to\infty]{} 0.$