# Efron-Stein and McDiarmid Inequalities

# SO FAR

- Markov, Chebyshev & Chernoff
- Subgaussian & Subexponential
- (Exponential) tail bounds on $X$ & $\sum_{i=1}^{n} \alpha_i X_i$

Want: Tail bounds on functions of $(X_1 \cdots X_n)$ iid

① Largest eigenvalue of a random matrix

② Max degree of a random graph

Suppose $\quad X_1 \cdots X_n$

$$Z = g(X_1 \cdots X_n)$$

$$Z_i = \mathbb{E}[Z \mid X_1 \cdots X_i]$$

$$= \int g(X_1 \cdots X_i, x_{i+1} \cdots x_n) \, f_{X_{i+1} \cdots X_n \mid X_1 \cdots X_i}(x_{i+1} \cdots x_n) \, dx_{i+1} \cdots dx_n$$

$$Z_0 = \mathbb{E} g(X_1 \cdots X_n) = \mathbb{E} Z$$

$$Z_n = Z = g(X_1 \cdots X_n)$$

$$\mathrm{Var}(g(X_1 \cdots X_n)) = \mathbb{E}[(Z_n - Z_0)^2]$$

$$Z_n - Z_0 = \sum_{i=1}^{n} Z_i - Z_{i-1}$$

$$Var(Z) = \mathbb{E}\left[ \left( \sum_{i=1}^{n} \underbrace{(Z_i - Z_{i-1})}_{\Delta_i} \right)^2 \right]$$

$$= \mathbb{E}\left[ \left( \sum_{i=1}^{n} \Delta_i \right)^2 \right]$$

$$= \mathbb{E}\left[ \sum_{i=1}^{n} \Delta_i^2 + 2 \sum_{i=1}^{n} \sum_{j=i+1}^{n} \Delta_i \Delta_j \right] \quad -\text{①}$$

$$\Delta_i = Z_i - Z_{i-1}$$

$$Z_i = \mathbb{E}\left[ Z \,|\, X_1 \cdots X_i \right] = \mathbb{E}_{X_{i+1} \cdots X_n}\left[ Z \right]$$

$$Z_{i-1} = \mathbb{E}\left[ Z \,|\, X_1 \cdots X_{i-1} \right] = \mathbb{E}_{X_i \cdots X_n}\left[ Z \right]$$

$$Z_{i-1} = \mathbb{E}_{X_i \cdots X_n}[Z] = \mathbb{E}_{X_i} \mathbb{E}_{X_{i+1} \cdots X_n}[Z]$$

$$= \mathbb{E}_{X_i} Z_i$$

$$\Delta_i = Z_i - \mathbb{E}_{X_i} Z_i \quad = \quad \mathbb{E}_{X_{i+1}^n} f(X^n) - \mathbb{E}_{X_i^n} f(X^n)$$

$$\mathbb{E}_{X^n} \Delta_i = 0$$

Consider $\quad j > i \qquad \mathbb{E}_{X^n}(\Delta_i \Delta_j)$

$$\mathbb{E}_{X^n} \left[ \left( \mathbb{E}_{X_{i+1}^n} Z - \mathbb{E}_{X_i^n} Z \right) \times \left( \mathbb{E}_{X_{j+1}^n} Z - \mathbb{E}_{X_j^n} Z \right) \right]$$

$$X_{i+1}^n = (X_{i+1} \cdots X_n)$$

$$\underbrace{\mathbb{E}_{X_i \cdot X_{j} \cdot X_{j}, X_i} \mathbb{E}_{X_j}}_{\mathbb{E}_{X^n}} \left[ \underbrace{\left( \mathbb{E}_{X_{i+1}^n} Z - \mathbb{E}_{X_i^n} Z \right)}_{\substack{\text{not a function} \\ \text{of } X_j \\ \text{for } j > i}} \times \left( \mathbb{E}_{X_{j+1}^n} Z - \mathbb{E}_{X_j^n} Z \right) \right]$$

$$= \left( \mathbb{E}_{X_{i+1}^n} Z - \mathbb{E}_{X_i^n} Z \right) \mathbb{E}_{X_j} \underbrace{\left( \mathbb{E}_{X_{j+1}^n} Z - \mathbb{E}_{X_j^n} Z \right)}_{\mathbb{E}_{X_j} \Delta_j}$$

$$\Rightarrow \mathbb{E}_{X^n} \Delta_i \Delta_j = 0 \quad \text{for } j > i$$

from ①,

$$\text{Var}(z) = \mathbb{E}\left[\sum_{i=1}^{n} \Delta_i^2\right]$$

$$= \sum_{i=1}^{n} \text{Var}(\Delta_i) \longrightarrow \text{true even}$$

for non iid $X_1 \cdots X_n$

$$z_1 \cdots z_n \qquad (X_1 \cdots X_n)$$

If $\mathbb{E}\left[z_{i+1} \mid X_1 \cdots X_i\right] = z_i \qquad \forall i,$

we say that $z_1 \cdots z_n$ is a martingale wrt

$$X_1 \cdots X_n$$

# Theorem (Efron-Stein inequality)

If $X_1 \cdots X_n$ independent & have bdd variance

$$g : \mathcal{X}^n \to \mathbb{R}$$

$$Z = g(X_1 \cdots X_n)$$

$$\mathrm{Var}(Z) \leq \sum_{i=1}^{n} \mathbb{E}\left(Z - \mathbb{E}_{X_i} Z\right)^2 = \sum_{i=1}^{n} \mathbb{E}_{X^n}\left(f(X^n) - \mathbb{E}_{X_i} f(X^n)\right)$$

$$\underbrace{\qquad\qquad}_{Y_i = Z - \mathbb{E}_{X_i} Z}$$

$$= \sum_{i=1}^{n} \mathrm{Var}(Y_i)$$

Define $Z_i' = g(X_1 \cdots X_{i-1}, X_i', X_{i+1} \cdots X_n)$

$\downarrow$

replace $X_i$ with an iid copy

$$\text{Var}(z) \leq \frac{1}{2} \sum_{i=1}^{n} \mathbb{E}(z - z_i')^2$$

$$\leq \frac{1}{2} \sum_{i=1}^{n} \mathbb{E}\left([z - z_i']_+\right)^2 \qquad \bigg| \begin{array}{l} [\alpha]_+ = \\ \quad \max\{\alpha, 0\} \end{array}$$

$$= \frac{1}{2} \sum_{i=1}^{n} \mathbb{E}\left([z - z_i']_-\right)^2 \qquad \bigg| \begin{array}{l} [\alpha]_- = \\ \quad -[-\alpha]_+ \end{array}$$

$$\leq \sum_{i=1}^{n} \inf_{\substack{U_i = h(x_1 \cdots x_{i-1}, x_{i+1} \cdots x_n) \\ \text{Var}(U_i) < \infty}} \mathbb{E}\left[(z - U_i)^2\right]$$

Suppose

$$|g(x_1 - x_n) - g(x_1 \cdots x_{i-1}, x_i', x_{i+1} \cdots x_n)| \le c_i$$

$$\text{Var}(z) = \text{Var}(g(x_n)) \le \frac{1}{2} \sum_{i=1}^{n} c_i^2$$

# Proof of Efron-Stein inequality

$$\text{Var}(Z) = \sum_{i=1}^{n} \mathbb{E}\Delta_i^2 \leq \sum_{i=1}^{n} \mathbb{E}_{X^n}\left( g(X^n) - \mathbb{E}_{X_i} g(X^n) \right)^2$$

$$\Delta_i = \mathbb{E}_{X_{i+1}^n} g(X^n) - \mathbb{E}_{X_i^n} g(X^n)$$

$$= \mathbb{E}_{X_{i+1}^n}\left[ g(X^n) - \mathbb{E}_{X_i} g(X^n) \right]$$

$$\Delta_i^2 = \left[ \mathbb{E}_{X_{i+1}^n}\left( g(X^n) - \mathbb{E}_{X_i} g(X^n) \right) \right]^2$$

$$\leq \mathbb{E}_{X_{i+1}^n}\left[ \underbrace{\left( g(X^n) - \mathbb{E}_{X_i} g(X^n) \right)^2}_{Y_i} \right]$$

$$\overline{\mathbb{E}_{X^n}} \Delta_i^2 \leq \mathbb{E}_{X^n}\left[ \left( g(X^n) - \mathbb{E}_{X_i} g(X^n) \right)^2 \right]$$

$$\mathbb{E} Y_i^2 \geq$$
$$\left( \mathbb{E} Y_i \right)^2$$
$$\Downarrow$$
$$\text{Var}(Y_i) \geq 0$$

$$\text{Var}\left(g(x^n)\right) \leq \sum_{i=1}^{n} \mathbb{E}_{X^n} \underbrace{\left(g(x^n) - \mathbb{E}_{X_i} g(x^n)\right)^2}$$

$$Z_i' = g(X_1 \cdots X_{i-1}, X_i', X_{i+1} \cdots X_n)$$

$\downarrow$ *iid copy*

<u>C1</u>: Conditioned on $X_1 \cdots X_{i-1}, X_{i+1} \cdots X_n$, $Z, Z_i'$ are iid

$$\mathbb{E}_{X_i X_i'} \left(Z_i - Z_i'\right)^2 = 2 \mathbb{E}_{X_i} \left(Z_i - (\mathbb{E}_{X_i} Z)\right)^2$$

Claim : $X \& Y$ iid

$$\text{Var}(X) = \tfrac{1}{2}\mathbb{E}(X - Y)^2$$

If $C_1$ is true,

$$\text{Var}(z) \leq \frac{1}{2} \sum_{i=1}^{n} \mathbb{E}_{x^n, x_i'} (z - z_i')^2$$

Proof of claim:

$$\frac{1}{2} \mathbb{E}(X - Y)^2$$

$$= \frac{1}{2} \mathbb{E}\left[(X - \mathbb{E}X) - (Y - \mathbb{E}Y)\right]^2$$

$$= \frac{1}{2} \left[ \mathbb{E}(X - \mathbb{E}X)^2 + \mathbb{E}(Y - \mathbb{E}Y)^2 - 2 \underbrace{\mathbb{E}\left[(X - \mathbb{E}X)(Y - \mathbb{E}Y)\right]}_{0} \right]$$

$$= \text{Var}(X)$$

HW! Show that if $X \& Y$ are iid,

$$\frac{1}{2} \mathbb{E}(X-Y)^2 = \mathbb{E}\left((X-Y)_+\right)^2$$

$$= \mathbb{E}\left((X-Y)_-\right)^2$$

$$(X-Y)_+ = \max\left\{(X-Y), 0\right\}$$

$$(X-Y)_- = \min\left\{(X-Y), 0\right\}$$

# Example: Bin packing problem

$$0.4, \quad 0.3, \quad 0.6, \quad 0.5$$

$$X_1 \cdots X_n \quad \text{iid} \qquad X_i \in [0, 1]$$

Goal: pack $X_1 \cdots X_n$ into min # bins.

(each bin has size $= 1$)

$$g(X_1 \sim X_n) = \text{min} \ \# \ \text{bins required}$$

Changing $x_i$ can change $g(x_1 \sim x_n)$ by at most $1$

$$\text{Var}(g(X_1 \sim X_n)) \leq \frac{n}{4}$$

**Example 2:** Longest common Subsequence

$$X^n = r \underline{e} \underline{c} \underline{e} \underline{n} t$$

$$Y^n = \underline{e}x\underline{c}e\underline{l}l\underline{e}\underline{n}t$$

$$g(X^n, Y^n) = \text{length of longest common Subsequence}$$

$X^n, Y^n$ iid

$$\frac{\mathbb{E} g(X^n, Y^n)}{n} \quad \text{is conjectured} \quad \frac{2}{1+\sqrt{2}} \quad \text{for } Ber\left(\frac{1}{2}\right)$$

$$Var(g(X^n, Y^n)) \leq \frac{n}{2}$$

$$P_n \left[ \; |g(x^n, y^n) - \mathbb{E}\, g(x^n, y^n)| > \delta \, \mathbb{E}\, g(x^n, y^n) \right]$$

$$\leq \quad \frac{\text{Var}(g)}{\delta^2 (\mathbb{E}\, g)^2} \quad \leq \quad \frac{\frac{n}{2}}{\delta^2 (Cn)^2}$$

$$= \quad \frac{1}{C\delta^2} \frac{1}{n}$$

# McDiarmid's inequality

If $X_1 \cdots X_n$ are independent

$$g: \mathcal{X}^n \to \mathbb{R}$$

$$| g(x_1 \cdots x_n) - g(x_1 \cdots x_{i-1}, x_i', x_{i+1} \cdots x_n) | \leq c_i$$

$$\forall x^n, x_i$$

then,

$$Pr\left[ |g(x^n) - \mathbb{E}g(x^n)| > t \right] \leq e^{-2t^2 / \sum_{i=1}^{n} c_i^2}$$

In fact, suppose

$$\sum_{i=1}^{n} (z - z_i')^2 \leq v^2 \qquad \text{with prob } 1,$$

$$Pr\left[ |g(x^n) - \mathbb{E}g(x^n)| > t \right] \leq e^{-t^2 / v^2}$$
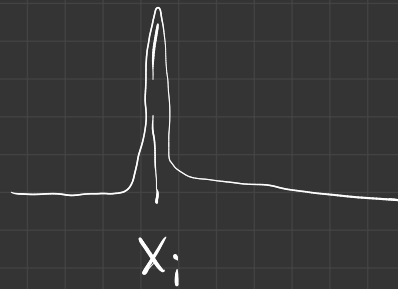
# Kernel density estimation

$$X_1 \cdots X_n \sim iid \ f_X$$

$\hookrightarrow$ unknown

$$K: \mathbb{R} \to \mathbb{R} \quad \text{smooth} \quad (\text{Kernel})$$

$$\phi_n(x) = \frac{1}{n h_n} \sum_{i=1}^{n} K\left(\frac{x - X_i}{h_n}\right)$$

$\underbrace{\phantom{xxxxxxxxx}}$

Mean $X_i$

Var $= h_n^2$

$$\int_{-\infty}^{\infty} K(x) \, dx = 1$$

$$K(x) \geq 0$$

$$\forall x$$

$X_i$

$$\text{Error} = \int_{-\infty}^{\infty} |f_X(x) - \phi_n(x)| \, dx$$

$$\Pr\left[\text{Error} > \mathbb{E}(\text{error}) + \delta\right]$$

$$g(x_1 - x_n) = \int_{-\infty}^{\infty} |f_X(x) - \phi_n(x)| \, dx$$

$$|g(x_1 - x_n) - g(x_1 - x_{i-1}, x_i', x_{i+1} - x_n)|$$

$$= \left| \int_{-\infty}^{\infty} |f_X(x) - \phi_n(x)| - |f_X(x) - \phi_n'(x)| \, dx \right|$$

$$\leq \int_{-\infty}^{\infty} |f_X(x) - \phi_n(x) - (f_X(x) - \phi_n'(x))| \, dx$$

$$= \int_{-\infty}^{\infty} |\phi_n'(x) - \phi_n(x)| \, dx$$

$$= \int_{-\infty}^{\infty} \left| \frac{1}{nh_n} \left( K\left(\frac{x-x_i}{h_n}\right) - K\left(\frac{x-x_i'}{h_n}\right) \right) \right| \, dx$$

$$\leq \int_{-\infty}^{\infty} \frac{1}{nh_n} \left( K\left(\frac{x-x_i}{h_n}\right) + K\left(\frac{x-x_i'}{h_n}\right) \right) \, dx$$

$$y = \frac{x-x_i}{h_n} \qquad dy = \frac{dx}{h_n}$$

$$\leq \int_{-\infty}^{\infty} \frac{1}{n} \left( K(y) + K(y') \right) \, dy = \frac{2}{n}$$

$$\text{Error} = \int_{-\infty}^{\infty} |f_X(x) - \phi_n(x)| \, dx \leq \int_{-\infty}^{\infty} (f_X(x) + \phi_n(x)) \, dx$$

$$\leq 2$$

$$\mathbb{E} \text{ error} \leq 2$$

$$\text{Var}(\text{error}) \leq \frac{1}{4} \times \sum_{i=1}^{n} c_i^2 \leq \frac{1}{n}$$

$$\text{Pr}\left[ \text{Error} \geq \overset{\frac{1}{\sqrt{n}}}{\mathbb{E} \text{ error}} (1+\delta) \right] \leq \frac{1}{\delta^2 n}$$

$$e^{-\delta^2 n}$$

# Empirical Risk Minimization

## Classification

$$X \qquad Y$$
$$\downarrow \qquad\qquad \downarrow$$
$$\text{image} \qquad \text{is object present}$$

$$X \in \mathcal{X} \qquad Y \in \{1, -1\}$$

$$(X, Y) \sim p_{XY}$$

Goal: Design $g: \mathcal{X} \to \{1, -1\}$
$$\downarrow$$
classifier

$$R_g = \Pr[g(X) \neq Y] \longrightarrow \text{risk for classifier } g$$

$$R_g = \mathbb{E} \, l(g(x), y)$$

$$\mathbb{E}(\text{risk}) = \mathbb{E} \, 1_{\{g(x) \neq y\}}$$

Suppose that we knew $p_{xy}$. What $g$ minimizes $R_g$?

$$g^*(x) = \underset{y \in \{1, -1\}}{\text{argmax}} \; p_{Y|x}(y \mid x) \quad \begin{pmatrix} MAP \\ \text{estimate} \end{pmatrix}$$

$$R(g^*) = \text{Minimum Bayes risk}$$

But we do not have $p_{xy}$

Dataset: $(X_1, Y_1) \, (X_2, Y_2) \sim (X_n, Y_n) \sim iid \, (p_{xy})$

$$R_n(g) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\{g(x_i) \neq y_i\}} \rightarrow \text{Empirical risk}$$

↓ empirical risk

$\underbrace{\phantom{\frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\{g(x_i) \neq y_i\}}}}$ fraction of time prediction is wrong

Empirical risk minimization

$$g_n = \underset{g \in \mathcal{Y}}{\arg\min} \; R_n(g)$$

What can we say about

$$|R(g) - R_n(g)| = \left| \mathbb{E}\, \mathbb{1}_{\{g(x) \neq y\}} - \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\{g(x_i) \neq y_i\}} \right|$$

$$\mathbb{E}\, R_n(g) = R(g)$$

$$R_n(g) = \frac{1}{n} \sum_{i=1}^{n} 1_{\{g(x_i) \neq y_i\}}$$

$$R(g) = \mathbb{E}_{X,Y} \, 1_{\{g(X) \neq Y\}} = Pr[g(X) \neq Y]$$

$$Pr\left[ |R_n(g) - R(g)| \geq \varepsilon R(g) \right] \leq 2 e^{-c \varepsilon^2 n R(g)}$$

Want :

$$R_n(g_n) - R(g^*)$$

optimum

$$\downarrow$$

$$\underset{g \in \mathcal{G}}{\text{argmin}} \, \frac{1}{n} \sum_{i=1}^{n} 1_{\{g(x_i) \neq y_i\}}$$

$$\text{argmin} \, \mathcal{Z}(g, x^n, y^n)$$

Toy example:    $(X, Y)$ -- $(X_n Y_n)$

$$y = \{g_1, g_{-1}\}$$

$$g_1(x) = 1 \qquad \forall x$$

$$g_{-1}(x) = -1 \qquad \forall x$$

$$P(Y = 1) = \alpha \qquad P(Y = -1) = 1 - \alpha$$

$$\sum_{i=1}^{n} 1_{\{g(x) \neq y_i\}} \begin{cases} Bin(n, 1-\alpha) & \text{if } g = g_1 \\ Bin(n, \alpha) & \text{if } g = g_{-1} \end{cases}$$

$$\min_{g \in \{g_1, g_{-1}\}} \sum_{i=1}^{n} 1_{\{g(x_i) \neq y_i\}} = \min\left\{ \begin{array}{l} \# \text{ of } 1's, \\ \#(-1)'s \end{array} \right.$$

$$\inf_{g \in \mathcal{G}} R_n(g)$$

$$R_n(g) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\{g(x_i) \neq y_i\}$$

$$g_n = \underset{g \in \mathcal{G}}{\arg\min} \; R_n(g)$$

$$g^* = \underset{g \in \mathcal{G}}{\arg\min} \; P_n[g(x) \neq y] \quad < \begin{array}{c} R(g^*) \\ R_n(g^*) \end{array}$$

$$R^* = \underset{g}{\arg\min} \; P_n[g(x) \neq y]$$

Observables:   $X^n, Y^n$   $g_n$   $R_n(g_n)$

Performance of $g_n$ on a (new) test sample:

$$P_n[g_n(x) \neq y] = R(g_n)$$

$$|R_n(g_n) - R(g_n)|$$

$$R(g_n) - R(g^*)$$

<span style="color:orange">min (test) risk for $g^=$ from ERM</span>

<span style="color:orange">(test) → min risk over all classifiers in $\mathcal{G}$</span>

Goal: ST $R(g_n) - R(g^*)$ is small

OR: What is min $n$ st $R(g_n) \approx R(g^*)$ whp

We know: for a given $g \in \mathcal{y}$

$$\Pr\left[\,|R_n(g) - R(g)| > \varepsilon R(g)\right] \leq e^{-cn\varepsilon^2}$$

for any $g$ 
$$R(g) = R_n(g) + R(g) - R_n(g)$$
$$\leq R_n(g) + \sup_{g \in \mathcal{y}} (R(g) - R_n(g))$$

Want:
$$\Pr\left[\sup_{g \in \mathcal{y}} (R(g) - R_n(g)) > \varepsilon\right]$$
$$\leq \sum_{g \in \mathcal{y}} \Pr\left[R(g) - R_n(g) > \varepsilon\right] \leq |\mathcal{y}|\, e^{-nc\varepsilon^2}$$

<u>Theorem</u>: If $\mathcal{y}$ is finite,
$$\Pr\left[R(g) \geq R_n(g) + 2\sqrt{\frac{\log|\mathcal{y}| + \log\frac{2}{\delta}}{2n}}\right] \leq \delta$$

$$\Pr\left[ R_n(g) \geq R(g^*) + 2\sqrt{\frac{\log|\mathcal{G}| + \log 2/\delta}{2n}} \right] \leq \delta$$

what if $\mathcal{Y}$ is infinite

$$\eta_g(x_i, y_i) = \mathbb{1}_{\{g(x_i) \neq y_i\}}$$

$$\mathcal{F}_{x^n, y^n} = \left\{ (\eta_g(x_1, y_1) \cdots \eta_g(x_n, y_n)) : g \in \mathcal{Y} \right\}$$

$$|\mathcal{F}_{x^n, y^n}| = 2^n$$

for given $\mathcal{Y}$

Growth function. $S_y = \sup_{(x^n, y^n)} \mathcal{F}_{x^n, y^n}$ → Measures how diverse $\mathcal{Y}$ is

**Theorem** (Vapnik -Cherronentis)

$$\Pr\left[ R(g) > R_n(g) + 2\sqrt{\frac{2 \log S_g(2n) + \log\frac{2}{\delta}}{n}} \quad \text{for any } g \in \mathcal{y}\right]$$

$$\leq \delta$$

" Introduction to Statistical Learning Theory "

Vapnik,    SLT