

Last time we proved Bennett's inequality: $\mathbb{E}X = 0$, $\mathbb{E}X^2 = \sigma^2$, $|X| < M = \text{const}$, X_1, \dots, X_n independent copies of X , and $t \geq 0$. Then

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{n\sigma^2}{M^2} \phi\left(\frac{tM}{n\sigma^2}\right)\right),$$

where $\phi(x) = (1+x)\log(1+x) - x$.

If X is small, $\phi(x) = (1+x)(x - \frac{x^2}{2} + \dots) - x = x + x^2 - \frac{x^2}{2} - x + \dots = \frac{x^2}{2} + \dots$.

If X is large, $\phi(x) \sim x \log x$.

We can weaken the bound by decreasing $\phi(x)$. Take¹ $\phi(x) = \frac{x^2}{2 + \frac{2}{3}x}$ to obtain **Bernstein's inequality**:

$$\begin{aligned} \mathbb{P}\left(\sum_{i=1}^n X_i \geq t\right) &\leq \exp\left(-\frac{n\sigma^2}{M^2} \left(\frac{(\frac{tM}{n\sigma^2})^2}{2 + \frac{2}{3}\frac{tM}{n\sigma^2}}\right)\right) \\ &= \exp\left(-\frac{t^2}{2n\sigma^2 + \frac{2}{3}tM}\right) \\ &= e^{-u} \end{aligned}$$

where $u = \frac{t^2}{2n\sigma^2 + \frac{2}{3}tM}$. Solve for t :

$$\begin{aligned} t^2 - \frac{2}{3}uMt - 2n\sigma^2u &= 0 \\ t &= \frac{1}{3}uM + \sqrt{\frac{u^2M^2}{9} + 2n\sigma^2u}. \end{aligned}$$

Substituting,

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq \sqrt{\frac{u^2M^2}{9} + 2n\sigma^2u} + \frac{uM}{3}\right) \leq e^{-u}$$

or

$$\mathbb{P}\left(\sum_{i=1}^n X_i \leq \sqrt{\frac{u^2M^2}{9} + 2n\sigma^2u} + \frac{uM}{3}\right) \geq 1 - e^{-u}$$

Using inequality $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$,

$$\mathbb{P}\left(\sum_{i=1}^n X_i \leq \sqrt{2n\sigma^2u} + \frac{2uM}{3}\right) \geq 1 - e^{-u}$$

For non-centered X_i , replace X_i with $X_i - \mathbb{E}X$ or $\mathbb{E}X - X_i$. Then $|X_i - \mathbb{E}X| \leq 2M$ and so with high probability

$$\sum (X_i - \mathbb{E}X) \leq \sqrt{2n\sigma^2u} + \frac{4uM}{3}.$$

Normalizing by n ,

$$\frac{1}{n} \sum X_i - \mathbb{E}X \leq \sqrt{\frac{2\sigma^2u}{n}} + \frac{4uM}{3n}$$

and

$$\mathbb{E}X - \frac{1}{n} \sum X_i \leq \sqrt{\frac{2\sigma^2u}{n}} + \frac{4uM}{3n}.$$

¹exercise: show that this is the best approximation

Whenever $\sqrt{\frac{2\sigma^2 u}{n}} \geq \frac{4uM}{3n}$, we have $u \leq \frac{n\sigma^2}{8M^2}$. So, $|\frac{1}{n} \sum X_i - \mathbb{E}X| \lesssim \sqrt{\frac{2\sigma^2 u}{n}}$ for $u \lesssim n\sigma^2$ (range of normal deviations). This is predicted by the Central Limit Theorem (condition for CLT is $n\sigma^2 \rightarrow \infty$). If $n\sigma^2$ does not go to infinity, we get Poisson behavior.

Recall from the last lecture that we're interested in concentration inequalities because we want to know $\mathbb{P}(f(X) \neq Y)$ while we only observe $\frac{1}{n} \sum_{i=1}^n I(f(X_i) \neq Y_i)$. In Bernstein's inequality take " X_i " to be $I(f(X_i) \neq Y_i)$. Then, since $2M = 1$, we get

$$\mathbb{E}I(f(X_i) \neq Y_i) - \frac{1}{n} \sum_{i=1}^n I(f(X_i) \neq Y_i) \leq \sqrt{\frac{2\mathbb{P}(f(X_i) \neq Y_i)(1 - \mathbb{P}(f(X_i) \neq Y_i))u}{n}} + \frac{2u}{3n}$$

because $\mathbb{E}I(f(X_i) \neq Y_i) = \mathbb{P}(f(X_i) \neq Y_i) = \mathbb{E}I^2$ and therefore $\text{Var}(I) = \sigma^2 = \mathbb{E}I^2 - (\mathbb{E}I)^2$. Thus,

$$\mathbb{P}(f(X_i) \neq Y_i) \leq \frac{1}{n} \sum_{i=1}^n I(f(X_i) \neq Y_i) + \sqrt{\frac{2\mathbb{P}(f(X_i) \neq Y_i)u}{n}} + \frac{2u}{3n}$$

with probability at least $1 - e^{-u}$. When the training error is zero,

$$\mathbb{P}(f(X_i) \neq Y_i) \leq \sqrt{\frac{2\mathbb{P}(f(X_i) \neq Y_i)u}{n}} + \frac{2u}{3n}.$$

If we forget about $2u/3n$ for a second, we obtain $\mathbb{P}(f(X_i) \neq Y_i)^2 \leq 2\mathbb{P}(f(X_i) \neq Y_i)u/n$ and hence

$$\mathbb{P}(f(X_i) \neq Y_i) \leq \frac{2u}{n}.$$

The above *zero-error rate* is better than $n^{-1/2}$ predicted by CLT.