# Lecture 7 Concerntration

IEMS 402 Statistical Learning

Northwestern

# Asymptotic VS Non-Asymptotic

# Drawback of Asymptotic Theory

Asymptotic : $r_n \left( T_n - \theta^* \right) \Rightarrow T$   $N[0, I_{\theta^*}^{-1})$

$\theta \in \mathbb{R}^d$ : $d$ is high.   Total Variance $\propto d$ $\Rightarrow$ $\sqrt{\dfrac{d}{n}}$

$\leftarrow$ Final Error

In most of the case.   $d \propto n$  . then  $\sqrt{\dfrac{d}{n}}$ is $O(1)$

$n \to \infty$

$O(1)$ distance to convergence.

# Concerntration

# First sense of Concerntration

inequalities of the form

Randomly sampled data #$n$

$$\mathbb{P}(X \geq t) \leq \phi(t) \qquad \phi(n, \epsilon)$$

Error/risk

where $\phi$ goes to zero (quickly) as $t \to \infty$

# First examples

Proposition (Markov's inequality)

If $X \geq 0$, then $\mathbb{P}(X \geq t) \leq \dfrac{\mathbb{E}[X]}{t}$ for all $t \geq 0$.

$$\mathbb{P}(X^2 \geq t^2) \leq \frac{\mathbb{E}[X^2]}{t^2}$$

(Markov's Inequality)

different convergence rate r.p.t  $t$

# First examples

Proposition (Markov's inequality)

If $X \geq 0$, then $\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[X]}{t}$ for all $t \geq 0$.

Proposition (Chebyshev's inequality)

For any $t \geq 0$, $\mathbb{P}(|X - \mathbb{E}[X]| \geq t) \leq \frac{\text{Var}(X)}{t^2}$

# First examples

Proposition (Markov's inequality)

If $X \geq 0$, then $\mathbb{P}(X \geq t) \leq \dfrac{\mathbb{E}[X]}{t}$ for all $t \geq 0$.

Proposition (Chebyshev's inequality)

For any $t \geq 0$, $\mathbb{P}(|X - \mathbb{E}[X]| \geq t) \leq \dfrac{\mathrm{Var}(X)}{t^2}$

$\downarrow$

$\dfrac{1}{t^3}$

Should be $O(e^{-t})$?

# Moment Generating Function

moment generating function

$$M_X(t) = \mathbb{E}[e^{tX}]$$

A function

Argument of moment generating function

$$e^x = 1 + t\,x + \frac{t^2 x^2}{2!} + \frac{t^3 x^3}{3!} + \dots$$

Control the reweighting the poly noirs.

# Moment Generating Function

moment generating function

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots$$

$$M_X(t) = \mathbb{E}[e^{tX}]$$

A function

Argument of moment generating function

$Y = X_1 + X_2$

$M_Y(t) = M_{X_1}(t) \, M_{X_2}(t)$

**Sum of Independent Random Variables:**

Suppose $X_1$, $X_2$, ..., $X_n$ are n independent random variables, and the random variable $Y$ is defined as

$$Y = X_1 + X_2 + \cdots + X_n.$$

Then,

$$
\begin{aligned}
M_Y(s) &= E[e^{sY}] \\
&= E[e^{s(X_1 + X_2 + \cdots + X_n)}] \\
&= E[e^{sX_1} e^{sX_2} \cdots e^{sX_n}] \\
&= E[e^{sX_1}]E[e^{sX_2}] \cdots E[e^{sX_n}] \quad \text{(since the } X_i\text{'s are independent)} \\
&= M_{X_1}(s)M_{X_2}(s) \cdots M_{X_n}(s).
\end{aligned}
$$

# Chernoff bound

$$P(X \geq a) \leq \inf_{t>0} M(t) e^{-ta}$$

$\downarrow$

$\mathbb{E}\left[ e^{tx} \right]$

$\Leftarrow$

$\mathbb{P}(X \geq a)$

$\|$

$\mathbb{P}\left( e^{tX} \geq e^{ta} \right) \leq \dfrac{\mathbb{E}\left[ e^{tx} \right]}{e^{ta}}$

$= M(t) \, e^{-ta}$

# Chernoff bound

$$P(X \geq a) \leq \inf_{t > 0} M(t) e^{-ta}$$

<u>reason 1</u>   $\sigma^2$ is "Variance"

## sub-Gaussian random variable

A mean-zero random variable $X$ is $\sigma^2$-sub-Gaussian if

$$\mathbb{E}[\exp(\lambda X)] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right) \quad \text{for all } \lambda \in \mathbb{R}.$$

<u>reason 2</u>

$X_1 + X_2 = Y$   $\rightarrow$ is $\sigma_1^2 + \sigma_2^2$ sub-Gaussian

$M_{X_1}(t) \, M_{X_2}(t) = M_Y(t)$

$\exp\left(\frac{t^2 \sigma_1^2}{2}\right)$   $\exp\left(\frac{t^2 \sigma_2^2}{2}\right)$

$\|$

$\exp\left(\frac{t^2 (\sigma_1^2 + \sigma_2^2)}{2}\right)$

Moment Generating function
of a Gaussian
Var is $\sigma^2$

Exercise

### Example
If $X \in [a, b]$, then

$$\mathbb{E}[\exp(\lambda(X - \mathbb{E}[X]))] \leq \exp\left(\frac{\lambda^2 (b-a)^2}{8}\right).$$

$\sigma = \frac{1}{4}(b-a)^2$

# Chernoff bound

$$P(X \geq a) \leq \inf_{t>0} M(t)e^{-ta}$$

sub-Gaussian random variable

A mean-zero random variable $X$ is $\sigma^2$-sub-Gaussian if

$$\mathbb{E}[\exp(\lambda X)] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right) \quad \text{for all } \lambda \in \mathbb{R}.$$

$$\mathbb{P}\left(X - \mathbb{E}[X] \geq t\right) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

# Hoeffding Inequality

$X_1 + \cdots, X_n$ is $\left(\sum_{i=1}^{n} \sigma_i^2\right) -$ sub gaussian.

## Corollary (Hoeffding bounds)

If $X_i$ are independent $\sigma_i^2$-sub-Gaussian random variables,

set probability to be $\Delta$me $O(1)$, $t = O\left(\frac{1}{\sqrt{n}}\right)$

$$\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^{n}(X_i - \mathbb{E}[X_i]) \geq t\right) \leq \exp\left(-\frac{nt^2}{\frac{2}{n}\sum_{i=1}^{n}\sigma_i^2}\right) \quad \rightarrow \text{this is a constant}$$

$$\| \quad$$

$$\mathbb{P}\left(\sum_{i=1}^{n}X_i - \mathbb{E}\left[\sum_{i=1}^{n}X_i\right] \geq nt\right) \leq \exp\left(-\frac{(nt)^2}{\sum_{i=1}^{n}\sigma_i^2}\right) = \exp\left(-\frac{nt^2}{\frac{2}{n}\sum_{i=1}^{n}\sigma_i^2}\right)$$

▶ usually stated as $X_i \in [a, b]$, so bound is $\exp\left(-\frac{2nt^2}{(b-a)^2}\right)$

> Should be $O(1/\sqrt{n})$?

# Moment Generating Function is Powerful

Bernstein's Inequality

$$\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^{n}X_i \geq t\right) \vee \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^{n}X_i \leq -t\right) \leq \exp\left(-\frac{nt^2}{2\sigma^2 + 2ct/3}\right),$$

different thing

$\sigma^2$ : variance     $|X_i| \leq c$

Special case: $\sigma$ is $0$

$\exp\left(-\frac{nt^2}{ct}\right)$

$= \exp\left(-\frac{nt}{c}\right) \rightarrow t=0 \left(\frac{1}{n}\right)$

**Homework 5, Question 3**

# Moment Generating Function is Powerful

**Proposition**

Let $\{Z_i\}_{i=1}^{N}$ be $\sigma^2$-sub-Gaussian (not necessarily independent). Then

$$\mathbb{E}\left[\max_i Z_i\right] \le \sqrt{2\sigma^2 \log N}.$$

$\Rightarrow$ max of $n$-random variables

is $\log N$ !!

$\exp\left(\sqrt{2\sigma^2 \log N}\right)$

$\downarrow$

$$\exp\left(\mathbb{E}\left[\max_i Z_i\right]\right) \le \mathbb{E}\left[\exp\left(t\max_i Z_i\right)\right]$$

$$\mathbb{E}\left[\sum_{i=1}^{n} \exp(t Z_i)\right] = O(N)$$

# Application

# Johnson-Lindenstrauss Lemma

**Lemma**  For any $0 < \epsilon < 1$ and any interger n let k be a possitive interger such that

$$k \geq \frac{24}{3\epsilon^2 - 2\epsilon^3} \log n \qquad (2)$$

error

$\epsilon$ is the error

number of projection

log of max of #(i,j)-pair r.v.

then for any set $A$ of $n$ points $\in \Re^d$ there exists a map $f : \Re^d \to \Re^k$ such that for all $x_i, x_j \in A$

$$(1 - \epsilon)||x_i - x_j||^2 \leq ||f(x_i) - f(x_j)||^2 \leq (1 + \epsilon)||x_i - x_j||^2 \qquad (3)$$

prob p fails of a (i,j) pair

$\Downarrow$

$\leq O(n^2 p)$ to fail the whole problem $\Rightarrow$ set p to be $O(\frac{1}{n^2})$

$||R(x_i - x_j)|| \approx ||x_i - x_j||$

How many (i,j) pair?   $O(n^2)$

https://cs.stanford.edu/people/mmahoney/cs369m/Lectures/lecture1.pdf
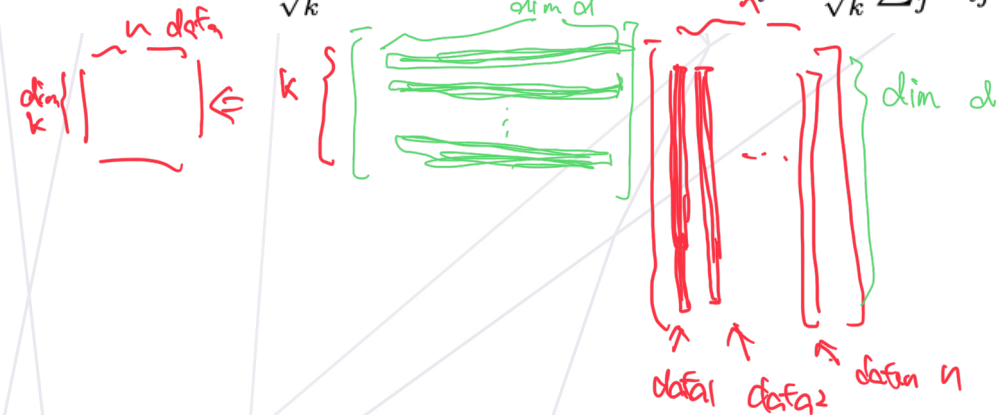
# Why it's important

CIFAR 100: 6000 32x32images,



f - linear

# Idea: random projection

**Definition** Let R be a random matrix of order $k \times d$ i.e $R_{ij} \overset{i.i.d}{\sim} N(0,1)$ and $u$ be any fixed vector $\in \Re^d$. Define $v = \frac{1}{\sqrt{k}} R \cdot u$. Thus $v \in \Re^k$ and $v_i = \frac{1}{\sqrt{k}} \sum_j R_{ij} u_j$

# Why it's important



SVD

Randomized SVD

Halko, Nathan, Per-Gunnar Martinsson, and Joel A. Tropp. "Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions." SIAM review
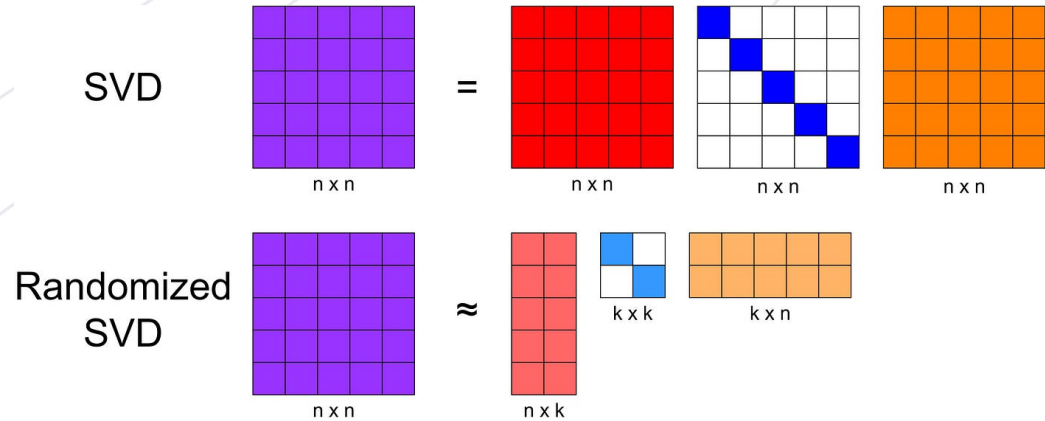
# Idea: random projection

**Definition** Let R be a random matrix of order $k \times d$ i.e $R_{ij} \overset{i.i.d}{\sim} N(0,1)$ and $u$ be any fixed vector $\in \Re^d$. Define $v = \frac{1}{\sqrt{k}} R \cdot u$. Thus $v \in \Re^k$ and $v_i = \frac{1}{\sqrt{k}} \sum_{j=1}^{k} R_{ij} u_j$

normalization

**Fact 1.** $\mathbb{E}[\|v\|^2] = \|u\|^2$ →

$u \in \mathbb{R}^d, \quad r \in \mathbb{R}^d \quad r \sim N(0,I)$

$\mathbb{E}\left[(u \cdot r)^2\right] = \mathbb{E}\left[u^T r r^T u\right] = u^T \underset{\text{covariance}}{\mathbb{E}\left[r r^T\right]} u = u^T I u = u^T u = \|u\|^2$

$\left[\sum_{i=1}^{d} u_i \cdot r_i\right]^2 \Rightarrow N(0, (u_1^2 + u_2^2 + \cdots + u_d^2)) \Rightarrow \mathbb{E}\left[(u \cdot r)^2\right] = \|u\|^2$

$N(0,1)$

**Question.** $\mathbb{P}(\|v\|^2 \geq (1 + \epsilon)\|u\|^2)$    Assume $\|u\| = 1$

# Random projection

**Question.** $\mathbb{P}(\|v\|^2 \geq (1 + \epsilon)\|u\|^2)$

Means $\dfrac{\sum_{i=1}^{k} x_i^2}{k} \geq (1 + \epsilon)$

$x_i = R_i^\top \cdot u$

$k$ of $d$ 1-dim projection

# Random projection

**Question.** $\mathbb{P}(\|v\|^2 \geq (1+\epsilon)\|u\|^2)$

$x_i = R_i^\top \cdot u$

Means $\dfrac{\sum_{i=1}^k x_i^2}{k} \geq (1+\epsilon) \rightarrow e^{\lambda x} \geq e^{\lambda(1+\epsilon)k}$

$\lambda = \sum_{i=1}^k x_i^2$

$\mathbb{E}[e^{\lambda x}] = \prod_{i=1}^k \mathbb{E}[e^{\lambda x_i^2}] = \left(\mathbb{E}[e^{\lambda x_i^2}]\right)^k$

# Random projection

**Question.** $\mathbb{P}(\|v\|^2 \geq (1+\epsilon)\|u\|^2)$

$$x_i = R_i^\top \cdot u$$

Means $\dfrac{\sum_{i=1}^{k} x_i^2}{k} \geq (1+\epsilon) \rightarrow e^{\lambda x} \geq e^{\lambda(1+\epsilon)k}$

$$\mathbb{E}[e^{\lambda x}] = \prod_{i=1}^{k} \mathbb{E}[e^{\lambda x_i^2}] = \left(\mathbb{E}[e^{\lambda x_i^2}]\right)^k$$

$\mathbb{E}[e^{\lambda x_i}]$ is the moment generating function of a $\chi^2$.

Thus $\mathbb{P}[e^{\lambda(1+\epsilon)k}] \leq \left(\dfrac{1}{\sqrt{1-2\lambda}}\right)^k \cdot \dfrac{1}{e^{\lambda(1+\epsilon)k}}$

# Random projection

**Question.** $\mathbb{P}(\|v\|^2 \geq (1+\epsilon)\|u\|^2) \leq e^{-(\epsilon^2/2 - \epsilon^3)k/2}$

$x_i = R_i^\top \cdot u$

Means $\dfrac{\sum_{i=1}^{k} x_i^2}{k} \geq (1+\epsilon) \rightarrow e^{\lambda x} \geq e^{\lambda(1+\epsilon)k}$

$$\mathbb{E}[e^{\lambda x}] = \prod_{i=1}^{k} \mathbb{E}[e^{\lambda x_i}] = \left(\mathbb{E}[e^{\lambda x_i}]\right)^k$$

Thus $\mathbb{P}[e^{\lambda(1+\epsilon)k}] \leq \left(\dfrac{1}{\sqrt{1-2\lambda}}\right)^k \cdot \dfrac{1}{e^{\lambda(1+\epsilon)k}}$

set $\lambda = \dfrac{\epsilon}{2(1+\epsilon)}$

$\leq e^{-(\epsilon^2/2 - \epsilon^3)k/2} \leq n^{-2}$

Why?

Uniform bound!

# Note

another proof using epsilon-net: **Theorem 8.**
https://www.cs.princeton.edu/~smattw/Teaching/Fa19Lectures/lec9/lec9.pdf