Probability Review/New Material High Dim Probability & Linear Algebra for ML and Sig Proc

Namrata Vaswani

Iowa State University

<ロト <部ト <注入 <注下 = 正

900

#### Reading material, Relevant Courses I

- The book "High Dimensional Probability for Data Science" by Roman Vershynin; and early parts
- 2 The tutorial article on "Non-asymptotic Random Matrix Theory" also by Vershynin
- Probability: https://www.ece.iastate.edu/~namrata/EE527\_Spring14/Probability\_recap\_3.pdf

Good courses to take at ISU: EE 523, STAT 542, 543, EE/Math 623X

Linear Algebra (based on first few chapters of Horn and Johnson, Matrix Analysis): https://www.ece.iastate.edu/~namrata/EE527\_Spring14/linearAlgebraNotes.pdf

Good courses to take at ISU: MATH 510 (first half of the course); if too advanced, then first take MATH 407/507.

#### 6 Review of Basics:

Probability:

https://www.ece.iastate.edu/~namrata/EE527\_Spring12/322\_recap.pdf http://cs229.stanford.edu/section/cs229-prob.pdf http://cs229.stanford.edu/section/more\_on\_gaussians.pdf

イロト 不得 トイラト イラト 二日

Linear Algebra: Andrew Ng's review from CS229 course at Stanford: http://cs229.stanford.edu/section/cs229-linalg.pdf also http://cs229.stanford.edu/livenotes2020spring/linearalgebra-slides.pdf

3

イロト 不得 トイヨト イヨト

Chapter 1 of book (Vershynin's book)

E

イロト イヨト イヨト イヨト

## Notation I

- Order etc
  - Order notation:  $f(n) \in O(g(n))$  means that there exists an  $n_0 < \infty$  such that for all  $n > n_0$ ,  $f(n) \le Cg(n)$  for a numerical constant C
  - Omega notation:  $f(n) \in \Omega(g(n))$  means that there exists an  $n_0 < \infty$  such that for all  $n > n_0$ ,  $f(n) \ge Cg(n)$  for a numerical constant C
  - $a \ll b$  means a/b is less than O(1)
  - Re-use of letter C: C is used to denote different numerical constants in different uses
- Linear algebra
  - For a matrix A, A' or A<sup>T</sup> or A<sup>T</sup> denotes matrix transpose; other use of MATLAB notation too.
  - Sphere in  $\Re^n$ :  $S^{n-1}$ , e.g., circle is a sphere in  $\Re^2$  and is denoted by  $S^1$
  - ▶ Norms: ||.||: I2-norm, ||.||<sub>1</sub>: I1-norm, ||.||<sub>F</sub>: Frobenius norm
  - ▶ Indicator function: 1 statement = 1 if statement is true and = 0 otherwise.
- Probability
  - ▶ For a set *A*, *A<sup>c</sup>* denotes its complement set.
  - Cumulative Distribution Function (CDF):  $F_X(x) := \Pr(X \le x)$
  - MGF  $M_X(t) = \mathbb{E}[e^{tX}]$  for a scalar X. For a vector,  $\underline{X}$ ,  $M_{\underline{X}}(\underline{u}) = \mathbb{E}[e^{\underline{u}'\underline{X}}]$
  - ► Characteristic function: C<sub>X</sub>(t) = E[e<sup>itX</sup>]: it is the FT of the distribution of X computed at frequency -t.

- $Pr(A, B) = Pr(A \text{ and } B) = Pr(A \cap B).$
- Gaussian *N*(μ, Σ)
- Bernoulli with probability of a 1 p: Bern(p)
- Symmetric Bernoulli SymBern: X = -1 w.p. 1/2 and X = +1 w.p. 1/2
- ▶ w.h.p. :
- ▶ w.p. :

3

<ロト <回ト < 回ト < 回ト < 回ト -

#### Basics: Simple algebra bounds: move to the end I

https://www.lkozma.net/inequalities\_cheat\_sheet/ineq.pdf

#### 2 Simple algebra bounds

- For any x > 0,  $1 + x < e^x$ used very often to convert  $\Pi_i(1 + \mu_i)$  to  $e^{\sum_i \mu_i}$  (appears when bounding MGF of sums of indep r.v.s)
- For 0 < x < 1,  $\log(1 + x) > x/(1 + x/2)$
- For all x > -1,  $\log(1 + x) \le x$
- For any x > 0,  $e^x < x + e^{x^2}$ used in subGaussian properties' equivalence. • ?? For any  $x \ge 0$ ,  $\frac{1}{1-x} \le e^{2x}$

For any 
$$z > 0$$
,  $\max(|z-1|, |z-1|^2) \le |z^2 - 1|$ 

- For any z > 0,  $|z 1| \ge \delta$  implies  $|z^2 1| \ge \max(\delta, \delta^2)$ .
- Stirling /factorial bounds

★ 
$$\Gamma(x) < x^x \Gamma(x) := ??$$

★  $p! > (p/e)^p$ , easy to see that  $p! < p^p$ 

★ 
$$(\frac{n}{k})^k \le {\binom{n}{k}} \le \sum_{k'=0}^k {\binom{n}{k}} \le (\frac{en}{k})^k$$

Taylor series

$$\star \exp(x) = \sum_{p=0}^{\infty} \frac{x^p}{p!}$$

Copy more from page 23, 30 of Vershynin book. TBD

◆ロト ◆帰 ト ◆臣 ト ◆臣 ト ◆ □ ●

#### Probability concepts assumed:

Probability axioms, disjoint events, independent events, conditional probability define, DeMorgan's laws, counting arguments

Use: try to convert an exact probability computation into probability of union of disjoint events, or intersection of independent events, or some combination of these ideas.

For upper bounding  $Pr(\cup_i A_i)$ : use union bound

For lower bounding  $Pr(\cup_i A_i)$ : use DeMorgan's + independence, and lower bounds on  $Pr(A_i)$  or use  $Pr(A) \ge Pr(A, B)$  followed by lower bound Pr(B) and Pr(A|B) (see use of this in the random vectors' theorem).

Many more ideas of course

Random variables: define PMF, joint PMF, PDF, joint PDF, CDF, joint CDF. Conditional CDF, conditional PDF.

Quick test of concepts: Given random variables (r.v.)  $X_1, X_2, \ldots X_n$ .

- **1** Compute distribution of  $Z = |X_1 + 1|$
- 2 Compute distribution of  $Z = X_1 \mod 5$  (remainder when  $X_1$  is divided by 5.
- 3 Compute the distribution of  $Z = X_1 + X_2$ . First
- Occupie the distribution of the smallest r.v.,  $Z = \min(X_1, X_2, ..., X_n)$ .
- Sompute the distribution of the second smallest r.v. (2nd order statistic).

#### **Probability Review I**

• Chain rule: extension of  $P(A_1 \cap A_2) = P(A_1)P(A_2|A_1)$ 

 $\Pr(A_1, A_2, \dots, A_n) = \Pr(A_1) \Pr(A_2 | A_1) \dots \Pr(A_k | A_1, A_2, \dots, A_k - 1) \dots \Pr(A_n | A_1, A_2, \dots, A_{n-1})$ 

**2** Total expectation theorem for events, Law of iterated expectations for r.v.s Consider events  $A_1, A_2, \ldots, A_n$  that form a partition of the sample space. Partition means: all the events are disjoint and their union forms the entire sample space. Simplest example of a partition is n = 2,  $A_1 = A$ ,  $A_2 = A^c$ . We have

$$\mathbb{E}[X] = \sum_{i} \mathbb{E}[X|A_i] \operatorname{Pr}(A_i)$$

If we set  $X = \mathbb{1}_E$  for an event  $E_{i}$ , the above gives the total probability result.

$$\Pr(E) = \sum_{i} \Pr(E|A_i) \Pr(A_i)$$

For two r.v.s X, Y (scalar or vector r.vs),

$$\mathbb{E}[g(X,Y)] = \mathbb{E}[\mathbb{E}[g(X,Y)|X]]$$

(here  $\mathbb{E}[.]$  takes expectation w.r.t. all r.v.s - here  $X, Y; \mathbb{E}[.|X]$  takes expectation conditioned on X.

Namrata Vaswani (Iowa State U.)



Independence and Conditional independence of events, r.v.s:

**1** Two events independent: Pr(A, B) = Pr(A)Pr(B)

**2** A set of *n* events is independent if for any subset  $S \subseteq [1, 2, ...n]$ ,

$$\Pr(\cap_{i\in S}A_i) = \prod_{i\in S}\Pr(A_i)$$

 $\bigcirc$  A set of *n* r.v.s,  $X_1, X_2, ..., X_n$  independent iff joint distribution is equal to product of marginals

$$F_{X_1,X_2,...,X_n}(x_1,x_2,...,x_n) = \prod_{i=1}^n F_{X_i}(x_i)$$

- **G** Conditional independence given Z = z for all  $z \in C$ : above holds conditioned on Z = z for all  $z \in C$ .
- **③** i.i.d. : independent and  $F_{X_i}(x) = F_{X_1}(x)$ , so that

$$F_{X_1,X_2,...X_n}(x_1,x_2,...x_n) = \prod_{i=1}^n F_{X_i}(x_i) = \prod_{i=1}^n F_{X_1}(x_i)$$

Namrata Vaswani (Iowa State U.)

**6** X, Y (scalars or vectors) independent implies

 $\mathbb{E}[g(X)h(Y)] = \mathbb{E}[g(X)]\mathbb{E}[h(Y)]$ 

- Conditionally independent given event C: above holds given event C. Same for conditional indep given a r.v.
- **3** X indep of  $Y, Z \Rightarrow X$  indep Y; and X conditionally indep Y given Z.
- Cauchy-Schwarz inequality
  - For two vectors  $v_1, v_2$ ,

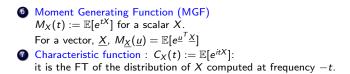
$$(v_1'v_2)^2 \le \|v_1\|^2 \|v_2\|^2$$

2 For two scalar r.v.s X, Y,

$$\mathbb{E}[XY]^2 \le \mathbb{E}[X^2]\mathbb{E}[Y^2]$$

- 3 obvious extensions apply for random vectors and matrices.
- **(5)** Union bound: for a set of events  $A_i$ , suppose that  $Pr(A_i) \ge 1 p_i$ . Then

$$\Pr(A_1, A_2, \dots, A_n) \equiv \Pr(\cap_i A_i) = 1 - \Pr(\cup_i A_i^c) \ge 1 - \sum_i \Pr(A_i^c) \ge 1 - \sum_i \exp(A_i^c) = \sum_i \exp(A_i^c) \ge 1 - \sum_i \exp(A_i^c) = \sum_i \exp(A_i^c) \ge 1 - \sum_i \exp(A_i^c) = \sum_i \exp(A_i^c) =$$



3

イロト 不得下 イヨト イヨト

# Scalar Gaussian r.v.

First note that a scalar Gaussian r.v. X with mean  $\mu$  and variance  $\sigma^2$  has the following pdf

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Its characteristic function can be computed by computing the Fourier transform at -t to get

$$C_X(t) = e^{j\mu t} e^{-\frac{\sigma^2 t^2}{2}}$$

# Gaussian random vector (Jointly Gaussian r.v.s)

Any of the following can be used as a definition of j G. All vectors should ideally be replaced by  $\underline{X}$  etc.

・ロト ・ 四 ト ・ 回 ト ・ 日 ト

The n × 1 random vector X is jointly Gaussian if and only if the scalar

#### $\mu^T X$

is Gaussian distributed for all  $n \times 1$  vectors u

2 The random vector X is jointly Gaussian if and only if its characteristic function,  $C_X(u) := \mathbb{E}[e^{iu^T X}]$  can be written as

$$C_X(u) = e^{iu^T \mu} e^{-u^T \Sigma u/2}$$

where  $\mu = \mathbb{E}[X]$  and  $\Sigma = cov(X)$ .

- ▶ Proof idea one side: Given X has above  $C_X(u)$ , show that  $V := u^T X$  is G for any vector u. To do this, show that  $C_V(t)$  has the G c.f. expression. To show this, use the fact that  $C_V(t) = C_X(tu)$  for scalar t.
- ▶ Proof idea other side: Given u'X is G for any u. Let V := u'X. Its mean and variance are  $\mu = u^T \mu$  and  $\sigma^2 = u^T \Sigma u$  and thus  $C_V(t) = e^{j\mu t} e^{-\frac{\sigma^2 t^2}{2}}$ . Now,  $C_X(u) = C_V(1) = e^{j\mu}e^{-\frac{\sigma^2}{2}}$ . Substituting for  $\mu, \sigma^2$  gives the  $C_X(u)$  expression we want to get.

The random vector X is j G if and only if it can be written as an affine function of i.i.d. standard Gaussian r.v's.

- Proof uses  $C_X(u)$  expression definition.
- ▶ Proof: suppose X = AZ + a where  $Z \sim \mathcal{N}(0, I)$ ; get an expression for its c.f. by using the c.f. definition and the fact that Z is a vector of i.i.d. standard Gaussian scalar r.v.s and thus  $\mathbb{E}[e^{itZ_j}] = e^{t^2/2}$  for any t. Show that the c.f. of X satisfies the  $C_X(u)$  formula given in 2 with  $\mu_X = a$ ,  $\Sigma_X = AA^T$ .
- Proof (other side): suppose X is j G with mean μ<sub>X</sub> and covariance Σ<sub>X</sub>; X can always be expressed as X = Σ<sup>1/2</sup>Z + μ where Z := Σ<sup>-1/2</sup>(X − μ); show that Z is std. G (by getting an expression for its c.f.).
   (c.f. of a std G Z is C<sub>Z</sub>(u) = e<sup>||u||<sup>2</sup>/2</sup>).
- The random vector X is j G if and only if it can be written as an affine function of jointly Gaussian r.v's.
  - Proof: Suppose X is an affine function of a j G r.v. Y, i.e. X = BY + b. Since Y is j G, by 3, it can be written as Y = AZ + a where Z ∼ N(0, I) (i.i.d. standard Gaussian). Thus, X = BAZ + (Ba + b), i.e. it is an affine function of Z, and thus, by 3, X is j G.
  - Proof (other side): X is j G. So by 3, it can be written as X = BZ + b. But  $Z \sim \mathcal{N}(0, I)$  i.e. Z is a j G r.v.

The random vector X is jointly Gaussian if and only if its joint pdf can be written as

$$f_X(x) = \frac{1}{(\sqrt{2\pi})^n det(\Sigma)} e^{-(X-\mu)^T \Sigma^{-1} (X-\mu)/2}$$
(1)

Proof: follows by computing the characteristic function from the pdf and vice versa. Suppose X has above PDF. Then  $C_X(u) = \mathbb{E}[\exp(iu'X)] = \int_x \exp(i\sum_j u_j x_j) f_X(x) dx$ , here x is a vector. Change of variables: let  $z = \Sigma^{-1/2}(x - \mu)$  and substitute into the integral. Integral will decouple into a product with term in the product being c.f. of a scalar Gaussian. Use formula, to finally get the vector Gaussian c.f. expression. Thus X is j G. Suppose X is j G. Then it has the given c.f. By uniqueness of Fourier transform, its density is given by (1).

# Properties

**1** If  $X_1, X_2$  are j G, then the conditional distribution of  $X_1$  given  $X_2$  is also j G

If the elements of a j G r.v. X are pairwise uncorrelated (i.e. non-diagonal elements of their covariance matrix are zero), then they are also mutually independent.

Any subset of X is also j G.

Integral identity

For a scalar r.v. Z that is non-negative, i.e.,  $Z \ge 0$  w.p. 1,

$$\mathbb{E}[Z] = \int_{\tau=0}^{\infty} \Pr(Z > \tau) d\tau$$

Proof: Use  $x = \int_{t=0}^{x} 1 dt = \int_{t=0}^{\infty} \mathbb{1}(t \le x) dt$  followed by moving expectation inside integral sign (allowed since indicator func is bounded).

**2** Use integral identity to convert w.h.p. bound to bound on expectation:

Given a non-negative r.v. Z that satisfies  $Pr(Z > u_0 + t) \le e^{-t^2}$  for all  $t \ge 0$  ( $Z \le 1.1u_0$  w.h.p.) for a  $u_0 \gg 2$  ( $u_0$  is more than order 1). This implies

 $\mathbb{E}[Z] \leq u_0 + 2 \leq 1.1 u_0$  the second bound assumes  $u_0 \gg 2$ 

- To use this second bound of  $1.1u_0$ , scale Z so that  $u_0 \gg 2$ .
- Similarly, if we are told that  $Pr(Z < u_0 t) \le e^{-t^2}$  for all  $t \ge 0$  ( $Z \ge 0.9u_0$  w.h.p.), assuming  $2 \ll u_0$ , we can show that

$$\mathbb{E}[Z] \ge u_0 - 2 \ge 0.9u_0$$

 Proof idea: apply integral identity, split integral into 0 to u<sub>0</sub> and then u<sub>0</sub> to ∞. In the first one, bound the probability by 1, in the second one, use the assumption, to get E[Z] ≤ u<sub>0</sub> + ∫<sub>t=0</sub><sup>∞</sup> e<sup>-t<sup>2</sup></sup> dt ≤ u<sub>0</sub> + √2π/2 < u<sub>0</sub> + 2. Proof idea for lower bound: E[Z] ≥ ∫<sub>u<sub>0</sub><sup>u<sub>0</sub></sub> Pr(Z > τ)dτ = ∫<sub>t=0</sub><sup>u<sub>0</sub></sup> Pr(Z > u<sub>0</sub> - t)dt ≥ ∫<sub>t=0</sub><sup>u<sub>0</sub></sup>(1 - e<sup>-t<sup>2</sup></sup>)dt
 Gaussian tail bounds: X ~ N(0, 1):
</sub></sup>

$$(\frac{1}{t} - \frac{1}{t^3})\frac{1}{\sqrt{2\pi}}e^{-t^2/2} \le \Pr(X \ge t) \le \frac{1}{\sqrt{2\pi}}\frac{1}{t}e^{-t^2/2}$$

Proof idea:

- ▶ Upper bound: use  $\int_{x=t}^{\infty} e^{-x^2/2} dx \leq \int_{x=t}^{\infty} (x/t) e^{-x^2/2} dx$  and then use change of variables to solve the integral.
- Lower bound: see page 12 of Vershynin book

For a non-negative r.v. Z,

$$\Pr(Z > s) \leq \frac{\mathbb{E}[Z]}{s}$$

Proof: easy application of integral identity

$$\mathbb{E}[Z] \geq \int_0^s \Pr(Z > \tau) d\tau \geq \Pr(Z > s) (\int_0^s d\tau) = \Pr(Z > s) s$$

#### Applications: basic ideas

- **(**) Apply this to  $Z = |X \mu|$  with  $\mu = \mathbb{E}[X]$ , to get Chebyshev's inequality.
- 2 Apply this to  $Z = e^{tX}$  for any  $t \ge 0$ . notice  $e^{tX}$  is always non-negative.

$$\Pr(X > s) = \Pr(e^{tX} > e^{ts}) \le e^{-ts} \mathbb{E}[e^{tX}] = e^{-ts} M_X(t)$$

Since this bound holds for all  $t \ge 0$ , we can take a min<sub>t \ge 0</sub> of the RHS or we can substitute in any convenient value of t.

3 To get a bound for 
$$Pr(X < -s)$$
, use  $Z = e^{-tX}$  for  $t \ge 0$ .

#### Markov inequality and applications II

**③** Useful for sums of independent r.v.s: if  $S = \sum_{i} X_{i}$  with  $X_{i}$ 's independent, then  $M_{X}(t) = \prod_{i} M_{X_{i}}(t)$ . So then we get

$$\Pr(\sum_{i} X_i > s) \le \min_{t \ge 0} e^{-ts} M_{\sum_i X_i}(t) = \min_{t \ge 0} e^{-ts} \prod_{i} \mathbb{E}[e^{tX_i}]$$

- Use exact expression for MGF or a bound on MGFs (e.g. Hoeffding's lemma bounds the MGF of any bounded r.v.)
- **(** Followed by often using  $1 + x \le e^x$  to simplify things
- **(**) Final step: either minimizer over  $t \ge 0$  by differentiating the expression or a pick a convenient value of t to substitute.
- **3** disregard this in first read: Final final step that is used sometimes: suppose get a bound g(s) but want to show  $g(s) \le f(s)$  for some simpler expression f(s): try to show that g(s) f(s) is a decreasing function for the desired range of s values with g(0) f(0) = 0 or something similar: this is used in Chernoff inequality for  $Bern(p_i)$  r.v.s. for small s setting.

#### Old recap document I

Quick test of concepts: Given random variables (r.v.)  $X_1, X_2, \ldots X_n$ .

- **1** Compute distribution of  $Z = |X_1 + 1|$
- 2 Compute distribution of  $Z = X_1 \mod 5$  (remainder when  $X_1$  is divided by 5.
- **3** Compute the distribution of  $Z = X_1 + X_2$ . First
- Occupate the distribution of the smallest r.v.,  $Z = \min(X_1, X_2, ..., X_n)$ .
- Ompute the distribution of the second smallest r.v. (2nd order statistic).

Some Topics

1 Chain Rule: extension of 
$$P(A_1 \cap A_2) = P(A_1)P(A_2|A_1)$$

 $P(A_1, A_2, \ldots, A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1, A_2) \ldots P(A_n|A_1, A_2, \ldots, A_{n-1})$ 

Total probability: if B<sub>1</sub>, B<sub>2</sub>,... B<sub>n</sub> form a partition of the sample space, then

$$P(A) = \sum_{i=1}^{n} P(A|B_i)P(B_i)$$

Partition: The events are mutually disjoint and their union is equal to the sample space.
Onion bound: suppose P(A<sub>i</sub>) ≥ 1 − p<sub>i</sub> for small probabilities p<sub>i</sub>, then

$$P(\cap_i A_i) = 1 - P(\cup_i A_i^c) \ge 1 - \sum_i P(A_i^c) \ge 1 - \sum_i p_i$$

Namrata Vaswani (Iowa State U.)

High Dim Prob & Lin Alg for ML

#### Old recap document II



Independence and Conditional Independence

events A. B are independent iff

$$P(A,B) = P(A)P(B)$$

• events  $A_1, A_2, \ldots, A_n$  are mutually independent iff for any subset  $S \subseteq \{1, 2, \ldots, n\}$ ,

$$P(\cap_{i\in S}A_i)=\prod_{i\in S}P(A_i)$$

- analogous definition for random variables: for mutually independent r.v.'s the joint pdf of any subset of r.v.'s is equal to the product of the marginal pdf's.
- events A, B are conditionally independent given an event C iff

$$P(A, B|C) = P(A|C)P(B|C)$$

- extend to a set of events as above
- extend to r v 's as above
- Side: Given X is independent of {Y, Z}. Then,
  - X is independent of Y; X is independent of Z

イロト 不得下 イヨト イヨト 二日

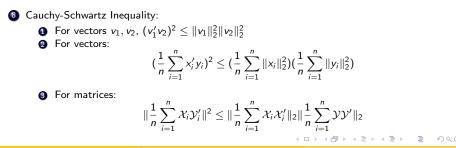
# Old recap document III

- X is conditionally independent of Y given Z
- $\blacktriangleright \mathbb{E}[XY|Z] = \mathbb{E}[X|Z]\mathbb{E}[Y|Z]$
- $\blacktriangleright \quad \mathbb{E}[XY|Z] = \mathbb{E}[X]\mathbb{E}[Y|Z]$
- 6 Law of Iterated Expectations:

$$\mathbb{E}_{X,Y}[g(X,Y)] = \mathbb{E}_{Y}[\mathbb{E}_{X|Y}[g(X,Y)|Y]]$$

Conditional Variance Identity:

$$Var_{X,Y}[g(X,Y)] = \mathbb{E}_{Y}[Var_{X|Y}[g(X,Y)|Y]] + Var_{Y}[\mathbb{E}_{X|Y}[g(X,Y)|Y]]$$



## Old recap document IV

- **(5)** For random vectors X, Y,

$$(\mathbb{E}[X'Y])^2 \le \mathbb{E}[||X||_2^2]\mathbb{E}[||Y||_2^2]$$

- **()** Proof follows by using the fact that  $\mathbb{E}[(X \alpha Y)^2] \ge 0$ . Get a quadratic equation in  $\alpha$  and use the condition to ensure that this is non-negative
- **7** For random matrices  $\mathcal{X}, \mathcal{Y}$ ,

$$\|\mathbb{E}[\mathcal{XY}']\|_2^2 \leq \lambda_{\sf max}(\mathbb{E}[\mathcal{XX}'])\lambda_{\sf max}(\mathbb{E}[\mathcal{YY}']) = \|\mathbb{E}[\mathcal{XX}']\|_2\|\mathbb{E}[\mathcal{YY}']\|_2$$

Recall that for a positive semi-definite matrix M,  $||M||_2 = \lambda_{\max}(M)$ . Proof: use the following definition of  $||M||_2$ :  $||M||_2 = \max_{x,y:||x||_2=1, ||y||_2=1} |x'My|$ , and then apply C-S for random vectors.

- Itoeffding's lemma: bounds the MGF of a zero mean and bounded r.v..
  - Suppose  $\mathbb{E}[X] = 0$  and  $P(X \in [a, b]) = 1$ , then

$$M_X(s):=\mathbb{E}[e^{sX}]\leq e^{rac{s^2(b-a)^2}{8}}$$
 if  $s>0$ 

Proof: use Jensen's inequality followed by mean value theorem, see http://www.cs.berkeley.edu/~jduchi/projects/probability\_bounds.pdf

#### Old recap document V

Oconvergence in probability. A sequence of random variables,  $X_1, X_2, ..., X_n$  converges to a constant *a* in probability means that for every  $\epsilon > 0$ ,

$$\lim_{n\to\infty}\Pr(|X_n-a|>\epsilon)=0$$

**4** Convergence in distribution. A sequence of random variables,  $X_1, X_2, \ldots, X_n$  converges to random variable Z in distribution means that

 $\lim_{n\to\infty} F_{X_n}(x) = F_Z(x), \text{ for almost all points} x$ 

Convergence in probability implies convergence in distribution

**(B)** Consistent Estimator. An estimator for  $\theta$  based on *n* random variables,  $\hat{\theta}_n(\underline{X})$ , is consistent if it converges to  $\theta$  in probability for large *n*.

Weak Law of Large Numbers (WLLN) for i.i.d. scalar random variables,  $X_1, X_2, \ldots X_n$ , with finite mean  $\mu$ . Define

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$$

For any  $\epsilon > 0$ ,

$$\lim_{n\to\infty} P(|\bar{X}_n-\mu|>\epsilon)=0$$

Proof: use Chebyshev if  $\sigma^2$  is finite. Else use characteristic function