EE 325: Probability and Random Processes Module 4: Bounds and Inequalities

D Manjunath

Networking Lab Department of Electrical Engineering Indian Institute of Technology

October 13, 2020

Topics in Module 4

- Recap of Union Bound and Bonferroni Inequalities
- Markov Inequality
- Chebyshev Inequality
- Chernoff Bound
- Cauchy-Schwarz Inequality
- Jensen's Inequality
- Hoeffding Bound

Recap Inequalities

- Often exact probabilities cannot be calculated and bounds may be more easily accessible. In fact, in many cases they suffice.
- Recall the **Union Bound:** If A_i , i = 1, ..., K are a set of K events, then

$$\Pr(\bigcup_{i=1}^K A_i) \le \sum_{i=1}^K \Pr(A_i)$$

- This is very useful when we know the individual probabilities but the union is hard to determine. Many a time, we may be just as happy with the upper bound from the Union Bound.
- From the Union Bound we can also write the following bound

$$\Pr(\cap_{i=1}^K A_i) \geq \sum_{i=1}^K (1 - \Pr(A_i))$$

 With a little more work, we can also obtain the lower following bound for the union of events.

$$\Pr(\bigcup_{i=1}^K A_i) \ge \sum_{i=1}^K \Pr(A_i) - \sum_{i < j} \Pr(A_i \cap A_j)$$

• This idea is useful: When you cannot obtain the exact probability.

An Identity

- Let X be discrete non negative random variable with pmf $p_X(i)$.
- Obviously, $p_X(i) = 0$ for i < 0.
- Let us evaluate $\sum_{i=0}^{\infty} \Pr(X > i)$

$$\sum_{i=0}^{\infty} \Pr(X > i) = p_X(1) + p_X(2) + \cdots$$

$$p_X(2) + p_X(3) + \cdots$$

$$p_X(3) + p_X(4) + \cdots$$

$$= \sum_{i=0}^{\infty} i p_X(i)$$

$$= E(X)$$

• This is also true for non negative continuous random variables

$$\mathsf{E}(\mathsf{X}) = \int_0^\infty (1 - F_\mathsf{X}(x)) \, dx$$



An Identity

- Let X be discrete non negative random variable with pmf $p_X(i)$.
- Obviously, $p_X(i) = 0$ for i < 0.
- Let us evaluate $\sum_{i=0}^{\infty} \Pr(X > i)$

$$\sum_{i=0}^{\infty} \Pr(X > i) = p_{X}(1) + p_{X}(2) + \cdots$$

$$p_{X}(2) + p_{X}(3) + \cdots$$

$$p_{X}(3) + p_{X}(4) + \cdots$$

$$= \sum_{i=0}^{\infty} i p_{X}(i)$$

$$= \exists (X)$$

• This is also true for non negative continuous random variables

$$\mathsf{E}(\mathsf{X}) = \int_0^\infty (1 - F_\mathsf{X}(x)) \, dx$$



An Identity

- Let X be discrete non negative random variable with pmf $p_X(i)$.
- Obviously, $p_X(i) = 0$ for i < 0.
- Let us evaluate $\sum_{i=0}^{\infty} \Pr(X > i)$

$$\sum_{i=0}^{\infty} \Pr(X > i) = p_X(1) + p_X(2) + \cdots$$

$$p_X(2) + p_X(3) + \cdots$$

$$p_X(3) + p_X(4) + \cdots$$

$$= \sum_{i=0}^{\infty} i p_X(i)$$

$$= E(X)$$

• This is also true for non negative continuous random variables

$$\mathsf{E}(\mathsf{X}) = \int_0^\infty (1 - F_\mathsf{X}(x)) \, dx$$



Markov Inequality

Clearly,

$$\mathsf{E}(\mathsf{X}) \ = \ \int (1 - F_{\mathsf{X}}(x)) \ dx \ \ge \ x \, \mathsf{Pr}(\mathsf{X} > x)$$

- This gives us the first key inequality.
- Markov Inequality: If X is a non negative random variable, then

$$\Pr(X > x) \le \frac{\mathsf{E}(X)}{x}$$

- Consider $X = (Y \mu_Y)^2$ where Y is a random variable.
- Clearly X is a non negative random variable. Hence we can apply Markov Inequality to X

$$\Pr(\mathsf{X} > x) \leq \frac{\mathsf{E}(X)}{x}$$

$$\Pr(\mathsf{Y} - \mu_{\mathsf{Y}})^2 > x) \leq \frac{\mathsf{E}(\mathsf{Y} - \mu_{\mathsf{Y}})^2}{x}$$

$$\Pr(\mathsf{Y} - \mu_{\mathsf{Y}}) > \sqrt{x}) \leq \frac{\mathsf{E}(\mathsf{Y} - \mu_{\mathsf{Y}})^2}{x}$$

$$\Pr(\mathsf{Y} - \mu_{\mathsf{Y}}) > y) \leq \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} \text{ (writing } y\sqrt{x})$$

$$\Pr(\mathsf{Y} - \mu_{\mathsf{Y}}) > y) \leq \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} = \frac{\sigma_{\mathsf{Y}}^2}{y^2}$$

- Consider $X = (Y \mu_Y)^2$ where Y is a random variable.
- Clearly X is a non negative random variable. Hence we can apply Markov Inequality to X

$$\Pr(\mathsf{Y} > x) \leq \frac{\mathsf{E}(X)}{x}$$

$$\Pr(\mathsf{Y} - \mu_{\mathsf{Y}})^2 > x) \leq \frac{\mathsf{E}(\mathsf{Y} - \mu_{\mathsf{Y}})^2}{x}$$

$$\Pr(\mathsf{Y} - \mu_{\mathsf{Y}} > \sqrt{x}) \leq \frac{\mathsf{E}(\mathsf{Y} - \mu_{\mathsf{Y}})^2}{x}$$

$$\Pr(\mathsf{Y} - \mu_{\mathsf{Y}} > y) \leq \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} \text{ (writing } y\sqrt{x})$$

$$\Pr(\mathsf{Y} - \mu_{\mathsf{Y}} > y) \leq \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} = \frac{\sigma_{\mathsf{Y}}^2}{y^2}$$

- Consider $X = (Y \mu_Y)^2$ where Y is a random variable.
- Clearly X is a non negative random variable. Hence we can apply Markov Inequality to X

$$\Pr(X > x) \leq \frac{E(X)}{x}$$

$$\Pr((Y - \mu_Y)^2 > x) \leq \frac{E((Y - \mu_Y)^2)}{x}$$

$$\Pr(|Y - \mu_Y| > \sqrt{x}) \leq \frac{E((Y - \mu_Y)^2)}{x}$$

$$\Pr(|Y - \mu_Y| > y) \leq \frac{VAR(Y)}{y^2} \text{ (writing } y\sqrt{x})$$

$$\Pr(|Y - \mu_Y| > y) \leq \frac{VAR(Y)}{y^2} = \frac{\sigma_Y^2}{y^2}$$

- Consider $X = (Y \mu_Y)^2$ where Y is a random variable.
- Clearly X is a non negative random variable. Hence we can apply Markov Inequality to X

$$\begin{aligned} & \text{Pr}(\mathsf{X} > x) & \leq & \frac{\mathsf{E}(X)}{x} \\ & \text{Pr}\Big((\mathsf{Y} - \mu_{\mathsf{Y}})^2 > x\Big) & \leq & \frac{\mathsf{E}\Big((\mathsf{Y} - \mu_{\mathsf{Y}})^2\Big)}{x} \\ & \text{Pr}\big(|\mathsf{Y} - \mu_{\mathsf{Y}}| > \sqrt{x}\big) & \leq & \frac{\mathsf{E}\Big((\mathsf{Y} - \mu_{\mathsf{Y}})^2\Big)}{x} \\ & \text{Pr}(|\mathsf{Y} - \mu_{\mathsf{Y}}| > y) & \leq & \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} \quad (\text{writing } y\sqrt{x}) \\ & \text{Pr}(|\mathsf{Y} - \mu_{\mathsf{Y}}| > y) & \leq & \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} & = \frac{\sigma_{\mathsf{Y}}^2}{y^2} \end{aligned}$$

- Consider $X = (Y \mu_Y)^2$ where Y is a random variable.
- Clearly X is a non negative random variable. Hence we can apply Markov Inequality to X

$$\begin{array}{rcl} \Pr(\mathsf{X} > x) & \leq & \frac{\mathsf{E}(X)}{x} \\ \Pr\left((\mathsf{Y} - \mu_{\mathsf{Y}})^2 > x\right) & \leq & \frac{\mathsf{E}\left((\mathsf{Y} - \mu_{\mathsf{Y}})^2\right)}{x} \\ \Pr\left(|\mathsf{Y} - \mu_{\mathsf{Y}}| > \sqrt{x}\right) & \leq & \frac{\mathsf{E}\left((\mathsf{Y} - \mu_{\mathsf{Y}})^2\right)}{x} \\ \Pr(|\mathsf{Y} - \mu_{\mathsf{Y}}| > y) & \leq & \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} \quad (\text{writing } y\sqrt{x}) \\ \Pr(|\mathsf{Y} - \mu_{\mathsf{Y}}| > y) & \leq & \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} & = \frac{\sigma_{\mathsf{Y}}^2}{y^2} \end{array}$$

- The event $\{|Y \mu_Y| > y\}$ is the same as $\{Y - \mu_Y < -v\} \cup \{Y - \mu_Y > v\}.$
- This in turn is the same as $\{Y\mu_Y < \mu_Y y\} \cup \{Y > \mu_Y + y\}$.
- Thus we have the Chebyshev Inequality

$$\begin{array}{lcl} \Pr(\mathsf{Y}\notin (\mu_{\mathsf{Y}}-y,\mu_{\mathsf{Y}}+y)) & \leq & \frac{\mathsf{VAR}(\mathsf{Y})}{y^2} \ = \ \frac{\sigma_{\mathsf{Y}}^2}{y^2} \\ \Pr(\mu_{\mathsf{Y}}-y \leq \mathsf{Y} \leq \mu_{\mathsf{Y}}+y)) & \geq & 1-\frac{\sigma_{\mathsf{Y}}^2}{y^2} \end{array}$$

EE 223

Chernoff Bound

- Consider $X = e^{sZ}$ where Z is a random variable. Assume s > 0,
- Clearly X is a non negative random variable. Hence we can apply Markov Inequality to X

$$\begin{array}{rcl} \Pr(\mathsf{X} > x) & \leq & \frac{\mathsf{E}(X)}{x} \\ \Pr(e^{s\mathsf{Z}} > x) & \leq & \frac{\mathsf{E}(e^{s\mathsf{Z}})}{x} \\ \Pr(e^{s\mathsf{Z}} > e^{s\mathsf{z}}) & \leq & \frac{\mathsf{E}(e^{s\mathsf{Z}})}{e^{s\mathsf{z}}} & (\text{writing } x = e^{s\mathsf{z}}) \end{array}$$

• This gives us the Chernoff Bound: For s > 0,

$$Pr(Z > z) \le \phi_X(s) e^{-sz}$$

• Note that this is true for all values of s! Thus we can obtain very tight bounds by choosing the value of s for which $\phi_X(s)$ e^{-sz} is minimum.



Chernoff Bound

- Consider $X = e^{sZ}$ where Z is a random variable. Assume s > 0,
- Clearly X is a non negative random variable. Hence we can apply Markov Inequality to X

$$\Pr(X > x) \leq \frac{E(X)}{x}$$

$$\Pr(e^{sZ} > x) \leq \frac{E(e^{sZ})}{x}$$

$$\Pr(e^{sZ} > e^{sz}) \leq \frac{E(e^{sZ})}{e^{sz}} \text{ (writing } x = e^{sz})$$

• This gives us the Chernoff Bound: For s > 0,

$$Pr(Z > z) \le \phi_X(s) e^{-sz}$$

• Note that this is true for all values of s! Thus we can obtain very tight bounds by choosing the value of s for which $\phi_X(s)$ e^{-sz} is minimum.



Chernoff Bound

- Consider $X = e^{sZ}$ where Z is a random variable. Assume s > 0,
- Clearly X is a non negative random variable. Hence we can apply Markov Inequality to X

$$\begin{array}{lcl} \Pr(\mathsf{X} > x) & \leq & \frac{\mathsf{E}(X)}{x} \\ \Pr(e^{s\mathsf{Z}} > x) & \leq & \frac{\mathsf{E}(e^{s\mathsf{Z}})}{x} \\ \Pr(e^{s\mathsf{Z}} > e^{sz}) & \leq & \frac{\mathsf{E}(e^{s\mathsf{Z}})}{e^{sz}} \quad (\text{writing } x = e^{sz}) \end{array}$$

• This gives us the Chernoff Bound: For s > 0,

$$Pr(Z > z) \le \phi_X(s) e^{-sz}$$

• Note that this is true for all values of s! Thus we can obtain very tight bounds by choosing the value of s for which $\phi_X(s)$ e^{-sz} is minimum.



• The vector calculus version of the Cauchy Schwarz Inequality is: If **u** and y are vectors, then

$$|u\cdot v|\leq |u|\;|v|$$

- Equality holds when the vectors collinear, i.e., $\mathbf{u} = a\mathbf{v}$. Furthermore,

$$f(a) = |au + v|^{2}$$

$$= \sum_{i} (au_{i} + v_{i})^{2} = \sum_{i} (a^{2}u_{i}^{2} + v_{i}^{2} + 2au_{i}v_{i})$$

$$= a^{2} \left(\sum_{i} u_{i}^{2}\right) + 2a \left(\sum_{i} u_{i}v_{i}\right) + \sum_{i} v_{i}^{2} \ge 0$$

$$4\left(\sum_{i}u_{i}v_{i}\right)^{2}-4\left(\sum_{i}u_{i}^{2}\right)\left(\sum_{i}u_{i}^{2}\right)\geq 0$$



• The vector calculus version of the Cauchy Schwarz Inequality is: If **u** and y are vectors, then

$$|u\cdot v|\leq |u|\;|v|$$

- Equality holds when the vectors collinear, i.e., $\mathbf{u} = a\mathbf{v}$. Furthermore, the LHS is zero if the vectors are orthhogonal.
- A proof is obtained by considering the following quadratic function in a.

$$f(a) = |au + v|^{2}$$

$$= \sum_{i} (au_{i} + v_{i})^{2} = \sum_{i} (a^{2}u_{i}^{2} + v_{i}^{2} + 2au_{i}v_{i})$$

$$= a^{2} \left(\sum_{i} u_{i}^{2}\right) + 2a \left(\sum_{i} u_{i}v_{i}\right) + \sum_{i} v_{i}^{2} \ge 0$$

$$4\left(\sum_{i}u_{i}v_{i}\right)^{2}-4\left(\sum_{i}u_{i}^{2}\right)\left(\sum_{i}u_{i}^{2}\right) \geq 0$$



The vector calculus version of the Cauchy Schwarz Inequality is: If u and y are vectors, then

$$|u\cdot v|\leq |u|\;|v|$$

- Equality holds when the vectors collinear, i.e., $\mathbf{u} = a\mathbf{v}$. Furthermore, the LHS is zero if the vectors are orthhogonal.
- A proof is obtained by considering the following quadratic function in

 a.

$$f(a) = |au + v|^{2}$$

$$= \sum_{i} (au_{i} + v_{i})^{2} = \sum_{i} (a^{2}u_{i}^{2} + v_{i}^{2} + 2au_{i}v_{i})$$

$$= a^{2} \left(\sum_{i} u_{i}^{2}\right) + 2a \left(\sum_{i} u_{i}v_{i}\right) + \sum_{i} v_{i}^{2} \ge 0$$

• Since this is non negative, the discriminant is non negative, i.e.

$$4\left(\sum_{i}u_{i}v_{i}\right)^{2}-4\left(\sum_{i}u_{i}^{2}\right)\left(\sum_{i}u_{i}^{2}\right)\geq 0$$

and the identity follows.



• The vector calculus version of the Cauchy Schwarz Inequality is: If **u** and y are vectors, then

$$|u\cdot v|\leq |u|\;|v|$$

- Equality holds when the vectors collinear, i.e., $\mathbf{u} = a\mathbf{v}$. Furthermore, the LHS is zero if the vectors are orthhogonal.
- A proof is obtained by considering the following quadratic function in a.

$$f(a) = |au + v|^{2}$$

$$= \sum_{i} (au_{i} + v_{i})^{2} = \sum_{i} (a^{2}u_{i}^{2} + v_{i}^{2} + 2au_{i}v_{i})$$

$$= a^{2} \left(\sum_{i} u_{i}^{2}\right) + 2a \left(\sum_{i} u_{i}v_{i}\right) + \sum_{i} v_{i}^{2} \ge 0$$

• Since this is non negative, the discriminant is non negative, i.e.,

$$4\left(\sum_{i}u_{i}v_{i}\right)^{2}-4\left(\sum_{i}u_{i}^{2}\right)\left(\sum_{i}u_{i}^{2}\right)\geq 0$$

and the identity follows.



• We have an analog of Cauchy Schwarz Inequality for random variables

$$\mathsf{E}((\mathsf{X} - \mu_{\mathsf{X}})(\mathsf{Y} - \mu_{\mathsf{Y}})) \le \sqrt{\mathsf{VAR}(X)}\,\mathsf{VAR}(Y)$$

For zero-mean random variables X and Y, this means

$$\mathsf{E}(\mathsf{X}\mathsf{Y}) \leq \sqrt{\mathsf{E}(\mathsf{X}^2)\,\mathsf{E}(\mathsf{Y}^2)}$$

- In fact we seen this before! Same as $-1 \le \rho_{XY} \le 1$
- A proof is exactly along the same lines as before

$$f(a) = \mathsf{E}\Big((a(\mathsf{X} - \mu_{\mathsf{X}}) + (\mathsf{Y} - \mu_{\mathsf{Y}}))^2\Big)$$

= $a^2\mathsf{E}\big((\mathsf{X} - \mu_{\mathsf{X}})^2\big) + \mathsf{E}\big((\mathsf{Y} - \mu_{\mathsf{Y}})^2\big) + 2\mathsf{E}((\mathsf{X} - \mu_{\mathsf{X}})(\mathsf{Y} - \mu_{\mathsf{Y}}))$
= $a^2\mathsf{VAR}(\mathsf{X}) + a(2\mathsf{COV}(\mathsf{X},\mathsf{Y})) + \mathsf{VAR}(\mathsf{Y})$

- f(a) is non negative because it is the expectation of non negative random variable and the final arguments are like before.
- Equality holds when Y = aX
- Using the analogy from vectors, we say that X and Y are orthogonal if COV(X, Y) = 0.



• We have an analog of Cauchy Schwarz Inequality for random variables

$$\mathsf{E}((\mathsf{X} - \mu_{\mathsf{X}})(\mathsf{Y} - \mu_{\mathsf{Y}})) \le \sqrt{\mathsf{VAR}(X)}\,\mathsf{VAR}(Y)$$

For zero-mean random variables X and Y, this means

$$\mathsf{E}(\mathsf{X}\mathsf{Y}) \leq \sqrt{\mathsf{E}(\mathsf{X}^2)\,\mathsf{E}(\mathsf{Y}^2)}$$

- In fact we seen this before! Same as $-1 \le \rho_{XY} \le 1$
- A proof is exactly along the same lines as before

$$f(a) = \mathsf{E}\Big((a(\mathsf{X} - \mu_{\mathsf{X}}) + (\mathsf{Y} - \mu_{\mathsf{Y}}))^2\Big)$$

= $a^2\mathsf{E}\big((\mathsf{X} - \mu_{\mathsf{X}})^2\big) + \mathsf{E}\big((\mathsf{Y} - \mu_{\mathsf{Y}})^2\big) + 2\mathsf{E}((\mathsf{X} - \mu_{\mathsf{X}})(\mathsf{Y} - \mu_{\mathsf{Y}}))$
= $a^2\mathsf{VAR}(\mathsf{X}) + a(2\mathsf{COV}(\mathsf{X},\mathsf{Y})) + \mathsf{VAR}(\mathsf{Y})$

- f(a) is non negative because it is the expectation of non negative random variable and the final arguments are like before.
- Equality holds when Y = aX
- Using the analogy from vectors, we say that X and Y are orthogonal if COV(X, Y) = 0.



• We have an analog of Cauchy Schwarz Inequality for random variables

$$\mathsf{E}((\mathsf{X} - \mu_{\mathsf{X}})(\mathsf{Y} - \mu_{\mathsf{Y}})) \le \sqrt{\mathsf{VAR}(X)}\,\mathsf{VAR}(Y)$$

For zero-mean random variables X and Y, this means

$$\mathsf{E}(\mathsf{X}\mathsf{Y}) \leq \sqrt{\mathsf{E}(\mathsf{X}^2)\,\mathsf{E}(\mathsf{Y}^2)}$$

- In fact we seen this before! Same as $-1 \le \rho_{XY} \le 1$
- A proof is exactly along the same lines as before

$$f(a) = \mathsf{E}\Big((a(\mathsf{X} - \mu_{\mathsf{X}}) + (\mathsf{Y} - \mu_{\mathsf{Y}}))^2\Big)$$

= $a^2\mathsf{E}\big((\mathsf{X} - \mu_{\mathsf{X}})^2\big) + \mathsf{E}\big((\mathsf{Y} - \mu_{\mathsf{Y}})^2\big) + 2\mathsf{E}((\mathsf{X} - \mu_{\mathsf{X}})(\mathsf{Y} - \mu_{\mathsf{Y}}))$
= $a^2\mathsf{VAR}(\mathsf{X}) + a(2\mathsf{COV}(\mathsf{X},\mathsf{Y})) + \mathsf{VAR}(\mathsf{Y})$

- f(a) is non negative because it is the expectation of non negative random variable and the final arguments are like before.
- Equality holds when Y = aX
- Using the analogy from vectors, we say that X and Y are orthogonal if COV(X,Y)=0.



Convex functions; Jensen's Inequality

• Recall that the expection of g(X) is

$$\mathsf{E}(g(\mathsf{X})) = \int g(x) f_{\mathsf{X}}(x) dx$$

• g(x) is a convex function if, for $0 \le a \le 1$,

$$g(ax_1 + (1-a)x_2) \le ag(x_1) + (1-a)g(x_2)$$

• Using induction, this can be generalized as follows. Let $a_1, a_2, \dots a_n$ be such that $a_i \ge 0$ for all i and $\sum_{i=1}^n a_i = 1$. If g(x) is convex, then

$$g(a_1x_1 + a_2x_2 + \dots + a_nx_n) \le a_1g(x_1) + a_2g(x_2) + \dots + a_ng(x_n)$$

- We can now use this as follows.
 - Assume X is a discrete random variable and it takes values $\{x_1, \dots, x_n\}$. Note the generalisation from X being integer valued.
 - Let $a_i = p_X(x_i)$. This gives us

$$g(p_{X}(x_{1})x_{1} + p_{X}(x_{2})x_{2} + \dots + p_{X}(x_{n})x_{n})$$

$$\leq p_{X}(x_{1})g(x_{1}) + p_{X}(x_{2})g(x_{2}) + \dots + p_{X}(x_{n})g(x_{n})$$

• LHS is g(E(X)) and RHS is E(g(X)). This gives **Jensen's Inequality**





Convex functions; Jensen's Inequality

• Recall that the expection of g(X) is

$$\mathsf{E}(g(\mathsf{X})) = \int g(x) f_{\mathsf{X}}(x) dx$$

• g(x) is a convex function if, for $0 \le a \le 1$,

$$g(ax_1 + (1-a)x_2) \le ag(x_1) + (1-a)g(x_2)$$

• Using induction, this can be generalized as follows. Let $a_1, a_2, \dots a_n$ be such that $a_i \ge 0$ for all i and $\sum_{i=1}^n a_i = 1$. If g(x) is convex, then

$$g(a_1x_1 + a_2x_2 + \dots + a_nx_n) \le a_1g(x_1) + a_2g(x_2) + \dots + a_ng(x_n)$$

- We can now use this as follows.
 - Assume X is a discrete random variable and it takes values $\{x_1, \dots, x_n\}$. Note the generalisation from X being integer valued.
 - Let $a_i = p_X(x_i)$. This gives us

$$g(p_{X}(x_{1})x_{1} + p_{X}(x_{2})x_{2} + \dots + p_{X}(x_{n})x_{n})$$

$$\leq p_{X}(x_{1})g(x_{1}) + p_{X}(x_{2})g(x_{2}) + \dots + p_{X}(x_{n})g(x_{n})$$

• LHS is g(E(X)) and RHS is E(g(X)). This gives **Jensen's Inequality**

$$g(\mathsf{E}(\mathsf{X})) \leq \mathsf{E}(g(\mathsf{X}))$$

DM EE 2

Jensen's Inequality

- There are several other proofs. Let us consider one more, which is also applicable for continuous random variables.
- Assume that g(x) is convex.
- Now consider the tangent to g(x) at E(X). Let mx + c be the equation describing this line.
 - Note that since this line passes through the point (E(X), g(E(X))), the following is true

$$g(\mathsf{E}(\mathsf{X})) = m\mathsf{E}(\mathsf{X}) + c.$$

- Convexity of g(x) means that the line obtained above is 'below' g(x), i.e., g(x) > mx + c for $a \le x \le b$.
- · We thus have

$$\mathsf{E}(g(\mathsf{X})) \ge \mathsf{E}(m\mathsf{X} + c) = m\mathsf{E}(\mathsf{X} + c) = g(\mathsf{E}(\mathsf{X}))$$



- Let X be a bounded zero-mean random variable, i.e., $\infty < a \le X \le b < \infty$ and E(X) = 0.
- Consider the function e^{sx} for $a \le x \le b$. This is a convex function. Hence

$$e^{sx} \le \frac{x-a}{b-a}e^{sb} + \frac{b-x}{b-a}e^{sa}$$

It helps to visualise this with a figure. Observe that $\frac{x-a}{b-a} + \frac{b-x}{b-a} = 1$ for all $a \le x \le b$.

$$E(e^{sX}) \leq E\left(\frac{X-a}{b-a}e^{sb}\right) + E\left(\frac{b-X}{b-a}e^{sa}\right)$$

$$= E\left(\frac{X}{b-a}e^{sb}\right) - E\left(\frac{a}{b-a}e^{sb}\right) + E\left(\frac{b}{b-a}e^{sa}\right) - E\left(\frac{X}{b-a}e^{sa}\right)$$

$$= E(X)\frac{1}{b-a}e^{sb} - \frac{a}{b-a}e^{sb} + \frac{b}{b-a}e^{sa} - E(X)\frac{1}{b-a}e^{sa}$$

$$E(e^{sX}) \leq \frac{be^{sa} - ae^{sb}}{ae^{sa}}$$

- Let X be a bounded zero-mean random variable, i.e., $\infty < a \le X \le b < \infty$ and E(X) = 0.
- Consider the function e^{sx} for $a \le x \le b$. This is a convex function. Hence

$$e^{sx} \le \frac{x-a}{b-a}e^{sb} + \frac{b-x}{b-a}e^{sa}$$

It helps to visualise this with a figure. Observe that $\frac{x-a}{b-a} + \frac{b-x}{b-a} = 1$ for all $a \le x \le b$.

$$E(e^{sX}) \leq E\left(\frac{X-a}{b-a}e^{sb}\right) + E\left(\frac{b-X}{b-a}e^{sa}\right)$$

$$= E\left(\frac{X}{b-a}e^{sb}\right) - E\left(\frac{a}{b-a}e^{sb}\right) + E\left(\frac{b}{b-a}e^{sa}\right) - E\left(\frac{X}{b-a}e^{sa}\right)$$

$$= E(X)\frac{1}{b-a}e^{sb} - \frac{a}{b-a}e^{sb} + \frac{b}{b-a}e^{sa} - E(X)\frac{1}{b-a}e^{sa}$$

$$E(e^{sX}) \leq \frac{be^{sa} - ae^{sb}}{b-a}$$

- Let X be a bounded zero-mean random variable, i.e., $\infty < a \le X \le b < \infty$ and E(X) = 0.
- Consider the function e^{sx} for $a \le x \le b$. This is a convex function. Hence

$$e^{sx} \le \frac{x-a}{b-a}e^{sb} + \frac{b-x}{b-a}e^{sa}$$

It helps to visualise this with a figure. Observe that $\frac{x-a}{b-a} + \frac{b-x}{b-a} = 1$ for all a < x < b.

$$E(e^{sX}) \leq E\left(\frac{X-a}{b-a}e^{sb}\right) + E\left(\frac{b-X}{b-a}e^{sa}\right)$$

$$= E\left(\frac{X}{b-a}e^{sb}\right) - E\left(\frac{a}{b-a}e^{sb}\right) + E\left(\frac{b}{b-a}e^{sa}\right) - E\left(\frac{X}{b-a}e^{sa}\right)$$

$$= E(X)\frac{1}{b-a}e^{sb} - \frac{a}{b-a}e^{sb} + \frac{b}{b-a}e^{sa} - E(X)\frac{1}{b-a}e^{sa}$$

$$E(e^{sX}) \leq \frac{be^{sa} - ae^{sb}}{ae^{sa}}$$

- Let X be a bounded zero-mean random variable, i.e., $\infty < a \le X \le b < \infty$ and E(X) = 0.
- Consider the function e^{sx} for $a \le x \le b$. This is a convex function. Hence

$$e^{sx} \le \frac{x-a}{b-a}e^{sb} + \frac{b-x}{b-a}e^{sa}$$

It helps to visualise this with a figure. Observe that $\frac{x-a}{b-a} + \frac{b-x}{b-a} = 1$ for all $a \le x \le b$.

$$E(e^{sX}) \leq E\left(\frac{X-a}{b-a}e^{sb}\right) + E\left(\frac{b-X}{b-a}e^{sa}\right)$$

$$= E\left(\frac{X}{b-a}e^{sb}\right) - E\left(\frac{a}{b-a}e^{sb}\right) + E\left(\frac{b}{b-a}e^{sa}\right) - E\left(\frac{X}{b-a}e^{sa}\right)$$

$$= E(X)\frac{1}{b-a}e^{sb} - \frac{a}{b-a}e^{sb} + \frac{b}{b-a}e^{sa} - E(X)\frac{1}{b-a}e^{sa}$$

$$E(e^{sX}) \leq \frac{be^{sa} - ae^{sb}}{b-a}$$

• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \le \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$f''(y) = \frac{x(1-x)e^{y}}{(x+(1-x)e^{y})^{2}} = \left(\frac{x}{(x+(1-x)e^{y})}\right) \left(1 - \frac{x}{(x+(1-x)e^{y})}\right)$$

• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \leq \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$f''(y) = \frac{x(1-x)e^{y}}{(x+(1-x)e^{y})^{2}} = \left(\frac{x}{(x+(1-x)e^{y})}\right) \left(1 - \frac{x}{(x+(1-x)e^{y})}\right)$$

• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \leq \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$f''(y) = \frac{x(1-x)e^{y}}{(x+(1-x)e^{y})^{2}} = \left(\frac{x}{(x+(1-x)e^{y})}\right)\left(1 - \frac{x}{(x+(1-x)e^{y})}\right)$$

• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \le \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$f''(y) = \frac{x(1-x)e^{y}}{(x+(1-x)e^{y})^{2}} = \left(\frac{x}{(x+(1-x)e^{y})}\right) \left(1 - \frac{x}{(x+(1-x)e^{y})}\right)$$

• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \leq \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$f''(y) = \frac{x(1-x)e^{y}}{(x+(1-x)e^{y})^{2}} = \left(\frac{x}{(x+(1-x)e^{y})}\right)\left(1 - \frac{x}{(x+(1-x)e^{y})}\right)$$

• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \leq \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

• Note f(0) = 0;. Performing Taylor series expansion around y = 0 we get, for some $z = \psi(y)$,

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$x(1-x)e^{y} \qquad (x-x)e^{y} \qquad ($$



DM

• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \leq \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

• Note f(0) = 0;. Performing Taylor series expansion around y = 0 we get, for some $z = \psi(y)$,

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$f''(y) = \frac{x(1-x)e^{y}}{(x+(1-x)e^{y})^{2}} = \left(\frac{x}{(x+(1-x)e^{y})}\right) \left(1 - \frac{x}{(x+(1-x)e^{y})}\right)$$

• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \le \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

• Note f(0) = 0;. Performing Taylor series expansion around y = 0 we get, for some $z = \psi(y)$,

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$f''(y) = \frac{x(1-x)e^{y}}{(x+(1-x)e^{y})^{2}} = \left(\frac{x}{(x+(1-x)e^{y})}\right)\left(1 - \frac{x}{(x+(1-x)e^{y})}\right)$$



• From the previous slide

$$\mathsf{E}\big(e^{s\mathsf{X}}\big) \le \frac{be^{sa} - ae^{sb}}{b - a} \ = \ xe^{sa} + (1 - x)e^{sb}$$

• Here we use x = b/(b-a). Let y = (b-a)s and consider

$$f(y) := \log (xe^{sa} + (1-x)e^{sb})$$

= $sa + \log (x + (1-x)e^{s(b-a)})$
= $(x-1)y + \log(x + (1-x)e^{y})$

• Note f(0) = 0;. Performing Taylor series expansion around y = 0 we get, for some $z = \psi(y)$,

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z(y))y^{2}$$

$$f'(y) = (x-1) + \frac{(1-x)e^{y}}{x+(1-x)e^{y}} \implies f'(0) = 0$$

$$f''(y) = \frac{x(1-x)e^{y}}{(x+(1-x)e^{y})^{2}} = \left(\frac{x}{(x+(1-x)e^{y})}\right)\left(1 - \frac{x}{(x+(1-x)e^{y})}\right)$$



Thus

$$f''(y) = \left(\frac{x}{(x+(1-x)e^y)}\right) \left(1 - \frac{x}{(x+(1-x)e^y)}\right) \le \frac{1}{4} \text{ for all } y.$$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z)y^2 = \frac{1}{2}f''(z)y^2$$

$$\le 0 + 0 + \frac{1}{2} \frac{1}{4} y^2 = \frac{y^2}{8}$$

$$\Xi(e^{sX}) \le e^{y^2/8} = e^{\frac{s^2(b-a)^2}{8}} \text{ recall that } y = (b-a)s$$

$$\mathsf{E}(e^{s\mathsf{X}}) \le e^{\frac{s^2(b-a)^2}{8}}$$



Thus

$$f''(y) = \left(\frac{x}{(x+(1-x)e^y)}\right) \left(1 - \frac{x}{(x+(1-x)e^y)}\right) \le \frac{1}{4} \text{ for all } y.$$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z)y^2 = \frac{1}{2}f''(z)y^2$$

$$\le 0 + 0 + \frac{1}{2} \frac{1}{4} y^2 = \frac{y^2}{8}$$

$$= (e^{x^2}) \le e^{x^2/8} = e^{\frac{x^2(b-a)^2}{2}} \text{ recall that } y = (b-a)s$$

$$\mathsf{E}(e^{s\mathsf{X}}) \le e^{\frac{s^2(b-a)^2}{8}}$$



Thus

$$f''(y) = \left(\frac{x}{(x+(1-x)e^y)}\right) \left(1 - \frac{x}{(x+(1-x)e^y)}\right) \le \frac{1}{4} \text{ for all } y.$$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z)y^2 = \frac{1}{2}f''(z)y^2$$

$$\le 0 + 0 + \frac{1}{2} \frac{1}{4} y^2 = \frac{y^2}{8}$$

$$= (e^{x^2}) \le e^{x^2/8} = e^{\frac{x^2(b-a)^2}{2}} \text{ recall that } y = (b-a)s$$

$$\mathsf{E}(e^{s\mathsf{X}}) \le e^{\frac{s^2(b-a)^2}{8}}$$



Thus

$$f''(y) = \left(\frac{x}{(x+(1-x)e^y)}\right) \left(1 - \frac{x}{(x+(1-x)e^y)}\right) \le \frac{1}{4} \text{ for all } y.$$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z)y^2 = \frac{1}{2}f''(z)y^2$$

$$\le 0 + 0 + \frac{1}{2} \frac{1}{4} y^2 = \frac{y^2}{8}$$

$$= (e^{xx}) \le e^{x^2/8} - e^{\frac{x^2}{2}(b-a)^2} \text{ recall that } y = (b-a)s$$

$$\mathsf{E}(e^{s\mathsf{X}}) \le e^{\frac{s^2(b-a)^2}{8}}$$



Thus

$$f''(y) = \left(\frac{x}{(x+(1-x)e^y)}\right) \left(1 - \frac{x}{(x+(1-x)e^y)}\right) \le \frac{1}{4} \text{ for all } y.$$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z)y^2 = \frac{1}{2}f''(z)y^2$$

$$\le 0 + 0 + \frac{1}{2}\frac{1}{4}y^2 = \frac{y^2}{8}$$

$$\mathsf{E}(e^{s\mathsf{X}}) \le e^{y^2/8} = e^{\frac{s^2(b-a)^2}{8}} \text{ recall that } y = (b-a)s$$

$$\mathsf{E}(e^{s\mathsf{X}}) \leq e^{\frac{s^2(b-a)^2}{8}}$$



Thus

$$f''(y) = \left(\frac{x}{(x+(1-x)e^y)}\right) \left(1 - \frac{x}{(x+(1-x)e^y)}\right) \le \frac{1}{4} \text{ for all } y.$$

$$f(y) = f(0) + f'(0)y + \frac{1}{2}f''(z)y^2 = \frac{1}{2}f''(z)y^2$$

$$\le 0 + 0 + \frac{1}{2}\frac{1}{4}y^2 = \frac{y^2}{8}$$

$$\mathsf{E}(e^{s\mathsf{X}}) \le e^{y^2/8} = e^{\frac{s^2(b-a)^2}{8}} \text{ recall that } y = (b-a)s$$

$$\mathsf{E}(e^{s\mathsf{X}}) \leq e^{\frac{s^2(b-a)^2}{8}}$$



- This is an extremely important inequality widely used in proving properties of machine learning algorithms.
- Consider X_i , i = 1, ..., n be independent bounded random variables with $a_i \le X_i \le b_i$. This means that $f_{X_i}(x) = 0$ for $x \in [a_i, b_i]$.
- We are interested in the tail distribution of $Y_n := \sum_{i=1}^n X_i$, i.e., $Pr(Y_n E(Y_n) > y)$. Use the Chernoff Bound with s > 0,

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) > y) \leq e^{-sy} \mathsf{E}\left(e^{s(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n))}\right)$$

$$= e^{-sy} \prod_{i=1}^n \mathsf{E}\left(e^{s(\mathsf{Y}_i - \mathsf{E}(\mathsf{Y}_i))}\right)$$

$$\leq e^{-sy} \prod_{i=1}^n e^{s^2(b_i - a_i)^2/8}$$

$$= e^{-sy + \frac{s^2}{8} \sum_{i=1}^n (b_i - a_i)^2}$$

- This is an extremely important inequality widely used in proving properties of machine learning algorithms.
- Consider X_i , i = 1, ..., n be independent bounded random variables with $a_i \le X_i \le b_i$. This means that $f_{X_i}(x) = 0$ for $x \in [a_i, b_i]$.
- We are interested in the tail distribution of $Y_n := \sum_{i=1}^n X_i$, i.e., $Pr(Y_n E(Y_n) > y)$. Use the Chernoff Bound with s > 0,

$$\begin{aligned} \mathsf{Pr}(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) > y) & \leq & e^{-sy} \mathsf{E} \Big(e^{s(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n))} \Big) \\ & = & e^{-sy} \prod_{i=1}^n \mathsf{E} \Big(e^{s(\mathsf{Y}_i - \mathsf{E}(\mathsf{Y}_i))} \Big) \\ & \leq & e^{-sy} \prod_{i=1}^n e^{s^2 (b_i - a_i)^2 / 8} \\ & = & e^{-sy + \frac{s^2}{8} \sum_{i=1}^n (b_i - a_i)^2} \end{aligned}$$



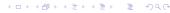
- This is an extremely important inequality widely used in proving properties of machine learning algorithms.
- Consider X_i , i = 1, ..., n be independent bounded random variables with $a_i \le X_i \le b_i$. This means that $f_{X_i}(x) = 0$ for $x \in [a_i, b_i]$.
- We are interested in the tail distribution of $Y_n := \sum_{i=1}^n X_i$, i.e., $Pr(Y_n E(Y_n) > y)$. Use the Chernoff Bound with s > 0,

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) > y) \leq e^{-sy} \mathsf{E}\left(e^{s(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n))}\right)$$

$$= e^{-sy} \prod_{i=1}^n \mathsf{E}\left(e^{s(\mathsf{Y}_i - \mathsf{E}(\mathsf{Y}_i))}\right)$$

$$\leq e^{-sy} \prod_{i=1}^n e^{s^2(b_i - a_i)^2/8}$$

$$= e^{-sy + \frac{s^2}{8} \sum_{i=1}^n (b_i - a_i)^2}$$



- This is an extremely important inequality widely used in proving properties of machine learning algorithms.
- Consider X_i , i = 1, ..., n be independent bounded random variables with $a_i \le X_i \le b_i$. This means that $f_{X_i}(x) = 0$ for $x \in [a_i, b_i]$.
- We are interested in the tail distribution of $Y_n := \sum_{i=1}^n X_i$, i.e., $Pr(Y_n E(Y_n) > y)$. Use the Chernoff Bound with s > 0,

$$\begin{aligned} \mathsf{Pr}(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) > y) & \leq & e^{-sy} \mathsf{E} \Big(e^{s(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n))} \Big) \\ & = & e^{-sy} \prod_{i=1}^n \mathsf{E} \Big(e^{s(\mathsf{Y}_i - \mathsf{E}(\mathsf{Y}_i))} \Big) \\ & \leq & e^{-sy} \prod_{i=1}^n e^{s^2(b_i - a_i)^2/8} \\ & = & e^{-sy + \frac{s^2}{8} \sum_{i=1}^n (b_i - a_i)^2} \end{aligned}$$



• Hoeffding Inequality: For bounded random variables X_i , for i = 1, ..., n, and $a_i \le X_i \le b_i$,

$$\Pr(\mathsf{Y}_{n} - \mathsf{E}(\mathsf{Y}_{n}) \ge y) \le e^{\frac{-2y^{2}}{\sum_{i=1}^{n} (b_{i} - a_{i})^{2}}}
\Pr(\mathsf{E}(\mathsf{Y}_{n}) - \mathsf{Y}_{n} \ge y) \le e^{\frac{-2y^{2}}{\sum_{i=1}^{n} (b_{i} - a_{i})^{2}}}$$

• If the X_i are iid, then

$$\Pr(Y_n - E(Y_n) \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

 $\Pr(E(Y_n) - Y_n \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

 $\Pr(\mathsf{E}(\mathsf{Y}_n) - \mathsf{Y}_n \ge y) \le e^{\frac{-2y^2}{n}}$



• Hoeffding Inequality: For bounded random variables X_i , for i = 1, ..., n, and $a_i \le X_i \le b_i$,

$$\begin{array}{lcl} \Pr(\mathsf{Y}_{n} - \mathsf{E}(\mathsf{Y}_{n}) \geq y) & \leq & e^{\frac{-2y^{2}}{\sum_{i=1}^{n}(b_{i}-a_{i})^{2}}} \\ \Pr(\mathsf{E}(\mathsf{Y}_{n}) - \mathsf{Y}_{n} \geq y) & \leq & e^{\frac{-2y^{2}}{\sum_{i=1}^{n}(b_{i}-a_{i})^{2}}} \end{array}$$

• If the X_i are iid, then

$$\Pr(Y_n - E(Y_n) \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

 $\Pr(E(Y_n) - Y_n \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y}{n}}$$

 $\Pr(\mathsf{E}(\mathsf{Y}_n) - \mathsf{Y}_n \ge y) \le e^{\frac{-2y}{n}}$



• Hoeffding Inequality: For bounded random variables X_i , for i = 1, ..., n, and $a_i \le X_i \le b_i$,

$$\Pr(\mathsf{Y}_{n} - \mathsf{E}(\mathsf{Y}_{n}) \ge y) \le e^{\frac{-2y^{2}}{\sum_{i=1}^{n} (b_{i} - a_{i})^{2}}}$$

$$\Pr(\mathsf{E}(\mathsf{Y}_{n}) - \mathsf{Y}_{n} \ge y) \le e^{\frac{-2y^{2}}{\sum_{i=1}^{n} (b_{i} - a_{i})^{2}}}$$

• If the X_i are iid, then

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

$$\Pr(\mathsf{E}(\mathsf{Y}_n) - \mathsf{Y}_n \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

 $\Pr(\mathsf{E}(\mathsf{Y}_n) - \mathsf{Y}_n \ge y) \le e^{\frac{-2y^2}{n}}$



• Hoeffding Inequality: For bounded random variables X_i , for i = 1, ..., n, and $a_i \le X_i \le b_i$,

$$\Pr(\mathsf{Y}_{n} - \mathsf{E}(\mathsf{Y}_{n}) \ge y) \le e^{\frac{-2y^{2}}{\sum_{i=1}^{n}(b_{i}-a_{i})^{2}}}$$

$$\Pr(\mathsf{E}(\mathsf{Y}_{n}) - \mathsf{Y}_{n} \ge y) \le e^{\frac{-2y^{2}}{\sum_{i=1}^{n}(b_{i}-a_{i})^{2}}}$$

• If the X_i are iid, then

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

$$\Pr(\mathsf{E}(\mathsf{Y}_n) - \mathsf{Y}_n \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

$$Pr(Y_n - E(Y_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

 $Pr(E(Y_n) - Y_n \ge y) \le e^{\frac{-2y^2}{n}}$



• Hoeffding Inequality: For bounded random variables X_i , for $i = 1, \ldots, n$, and $a_i < X_i < b_i$,

$$\Pr(\mathsf{Y}_{n} - \mathsf{E}(\mathsf{Y}_{n}) \ge y) \le e^{\frac{-2y^{2}}{\sum_{i=1}^{n} (b_{i} - a_{i})^{2}}}$$

$$\Pr(\mathsf{E}(\mathsf{Y}_{n}) - \mathsf{Y}_{n} \ge y) \le e^{\frac{-2y^{2}}{\sum_{i=1}^{n} (b_{i} - a_{i})^{2}}}$$

• If the X_i are iid, then

$$\Pr(Y_n - E(Y_n) \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

 $\Pr(E(Y_n) - Y_n \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

 $\Pr(\mathsf{E}(\mathsf{Y}_n) - \mathsf{Y}_n \ge y) \le e^{\frac{-2y^2}{n}}$



• Hoeffding Inequality: For bounded random variables X_i , for $i = 1, \ldots, n$, and $a_i < X_i < b_i$,

$$\begin{array}{lcl} \Pr(\mathsf{Y}_{n} - \mathsf{E}(\mathsf{Y}_{n}) \geq y) & \leq & e^{\frac{-2y^{2}}{\sum_{i=1}^{n}(b_{i} - a_{i})^{2}}} \\ \Pr(\mathsf{E}(\mathsf{Y}_{n}) - \mathsf{Y}_{n} \geq y) & \leq & e^{\frac{-2y^{2}}{\sum_{i=1}^{n}(b_{i} - a_{i})^{2}}} \end{array}$$

• If the X_i are iid, then

$$\Pr(Y_n - E(Y_n) \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

 $\Pr(E(Y_n) - Y_n \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

 $\Pr(\mathsf{E}(\mathsf{Y}_n) - \mathsf{Y}_n \ge y) \le e^{\frac{-2y^2}{n}}$



• Hoeffding Inequality: For bounded random variables X_i , for $i = 1, \ldots, n$, and $a_i < X_i < b_i$,

$$\begin{array}{lcl} \Pr(\mathsf{Y}_{n} - \mathsf{E}(\mathsf{Y}_{n}) \geq y) & \leq & e^{\frac{-2y^{2}}{\sum_{i=1}^{n}(b_{i} - a_{i})^{2}}} \\ \Pr(\mathsf{E}(\mathsf{Y}_{n}) - \mathsf{Y}_{n} \geq y) & \leq & e^{\frac{-2y^{2}}{\sum_{i=1}^{n}(b_{i} - a_{i})^{2}}} \end{array}$$

• If the X_i are iid, then

$$\Pr(Y_n - E(Y_n) \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$$

 $\Pr(E(Y_n) - Y_n \ge y) \le e^{\frac{-2y^2}{n(b-a)^2}}$

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

 $\Pr(\mathsf{E}(\mathsf{Y}_n) - \mathsf{Y}_n \ge y) \le e^{\frac{-2y^2}{n}}$



- A coin has an unknown bias p.
- The coin has been tossed n, let Y_n be the number of times a head has been observed.
- Note that $E(Y_n) = np$ and $E(Y_n/n) = p$.
- Hoeffding's Inequality tells us that

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

$$\Pr(\mathsf{Y}_n - np \ge ny) \le e^{-2ny^2}$$

$$\Pr\left(p - \frac{\mathsf{Y}_n}{n} \ge ny\right) \le e^{-2ny^2}$$

$$\Pr\left(p - \frac{\mathsf{Y}_n}{n} \ge y\right) \le e^{-2ny^2}$$

- Let us now see how to use this result in a practical setting.
- For example, let $e^{\frac{-2y^2}{n}} = 0.01$. This corresponds to $y \approx 2.14n$
- What can we say about the true value of *p* based on this observation?



- A coin has an unknown bias p.
- The coin has been tossed n, let Y_n be the number of times a head has been observed.
- Note that $E(Y_n) = np$ and $E(Y_n/n) = p$.
- Hoeffding's Inequality tells us that

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

$$\Pr(\mathsf{Y}_n - np \ge ny) \le e^{-2ny^2}$$

$$\Pr\left(p - \frac{\mathsf{Y}_n}{n} \ge ny\right) \le e^{-2ny^2}$$

$$\Pr\left(p - \frac{\mathsf{Y}_n}{n} \ge y\right) \le e^{-2ny^2}$$

- Let us now see how to use this result in a practical setting.
- For example, let $e^{\frac{-2y^2}{n}} = 0.01$. This corresponds to $y \approx 2.14n$
- What can we say about the true value of *p* based on this observation?



- A coin has an unknown bias p.
- The coin has been tossed n, let Y_n be the number of times a head has been observed.
- Note that $E(Y_n) = np$ and $E(Y_n/n) = p$.
- Hoeffding's Inequality tells us that

$$\Pr(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \ge y) \le e^{\frac{-2y^2}{n}}$$

$$\Pr(\mathsf{Y}_n - np \ge ny) \le e^{-2ny^2}$$

$$\Pr\left(p - \frac{\mathsf{Y}_n}{n} \ge ny\right) \le e^{-2ny^2}$$

$$\Pr\left(p - \frac{\mathsf{Y}_n}{n} \ge y\right) \le e^{-2ny^2}$$

- Let us now see how to use this result in a practical setting.
- For example, let $e^{\frac{-2y^2}{n}} = 0.01$. This corresponds to $y \approx 2.14n$
- What can we say about the true value of *p* based on this observation?



- A coin has an unknown bias *p*.
- The coin has been tossed n, let Y_n be the number of times a head has been observed.
- Note that $E(Y_n) = np$ and $E(Y_n/n) = p$.
- Hoeffding's Inequality tells us that

$$\begin{aligned} & \mathsf{Pr}(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \geq y) & \leq e^{\frac{-2y^2}{n}} \\ & \mathsf{Pr}(\mathsf{Y}_n - np \geq ny) & \leq e^{-2ny^2} \\ & \mathsf{Pr}\bigg(p - \frac{\mathsf{Y}_n}{n} \geq ny\bigg) & \leq e^{-2ny^2} \\ & \mathsf{Pr}\bigg(p - \frac{\mathsf{Y}_n}{n} \geq y\bigg) & \leq e^{-2ny^2} \end{aligned}$$

- This bounds the probability.
- Let us now see how to use this result in a practical setting.
- For example, let $e^{\frac{-2y^2}{n}} = 0.01$. This corresponds to $y \approx 2.14n$
- What can we say about the true value of p based on this observation?



- A coin has an unknown bias *p*.
- The coin has been tossed n, let Y_n be the number of times a head has been observed.
- Note that $E(Y_n) = np$ and $E(Y_n/n) = p$.
- Hoeffding's Inequality tells us that

$$\begin{aligned} & \mathsf{Pr}(\mathsf{Y}_n - \mathsf{E}(\mathsf{Y}_n) \geq y) & \leq e^{\frac{-2y^2}{n}} \\ & \mathsf{Pr}(\mathsf{Y}_n - np \geq ny) & \leq e^{-2ny^2} \\ & \mathsf{Pr}\bigg(p - \frac{\mathsf{Y}_n}{n} \geq ny\bigg) & \leq e^{-2ny^2} \\ & \mathsf{Pr}\bigg(p - \frac{\mathsf{Y}_n}{n} \geq y\bigg) & \leq e^{-2ny^2} \end{aligned}$$

- This bounds the probability.
- Let us now see how to use this result in a practical setting.
- For example, let $e^{\frac{-2y^2}{n}} = 0.01$. This corresponds to $y \approx 2.14n$
- What can we say about the true value of p based on this observation?



A learning problem

- There are three coins, named A, B, and C, each with unknownb biases possibly different. You are allowed a total of N tosses. For each toss you can choose any one of the three coins using any algorithm. And when you toss the chosen coin, you get a reward of one unit if it comes up heads.
- You have to 'learn' which of them is the more profitable coin and use it maximally.
- Before the n-th toss, let n_A , n_B , and n_C , be the number of times coins A, B, and C, respectively, have been tossed and let k_A , k_B , and k_C , be the number of heads for these coins.
- The preceding slide allows us to claim that value of p_A is less than $UCB_A = k_A/n_A + X_A$ with probability 0.01. X from n_A and the discussion in the previous slide. Similarly, for coins B and C.
- For the *n*-toss choosing the coin with the highest *UCB* has some very nice properties. We will investigate these in a computational experiment.