

# 18.175: Lecture 13

## More large deviations

Scott Sheffield

MIT

Legendre transform

Large deviations

Legendre transform

Large deviations

# Legendre transform

- ▶ Define **Legendre transform** (or Legendre dual) of a function  $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

# Legendre transform

- ▶ Define **Legendre transform** (or Legendre dual) of a function  $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

- ▶ Let's describe the Legendre dual geometrically if  $d = 1$ :  $\Lambda^*(x)$  is where tangent line to  $\Lambda$  of slope  $x$  intersects the real axis. We can “roll” this tangent line around the convex hull of the graph of  $\Lambda$ , to get all  $\Lambda^*$  values.

# Legendre transform

- ▶ Define **Legendre transform** (or Legendre dual) of a function  $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

- ▶ Let's describe the Legendre dual geometrically if  $d = 1$ :  $\Lambda^*(x)$  is where tangent line to  $\Lambda$  of slope  $x$  intersects the real axis. We can “roll” this tangent line around the convex hull of the graph of  $\Lambda$ , to get all  $\Lambda^*$  values.
- ▶ Is the Legendre dual always convex?

# Legendre transform

- ▶ Define **Legendre transform** (or Legendre dual) of a function  $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

- ▶ Let's describe the Legendre dual geometrically if  $d = 1$ :  $\Lambda^*(x)$  is where tangent line to  $\Lambda$  of slope  $x$  intersects the real axis. We can “roll” this tangent line around the convex hull of the graph of  $\Lambda$ , to get all  $\Lambda^*$  values.
- ▶ Is the Legendre dual always convex?
- ▶ What is the Legendre dual of  $x^2$ ? Of the function equal to 0 at 0 and  $\infty$  everywhere else?

# Legendre transform

- ▶ Define **Legendre transform** (or Legendre dual) of a function  $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

- ▶ Let's describe the Legendre dual geometrically if  $d = 1$ :  $\Lambda^*(x)$  is where tangent line to  $\Lambda$  of slope  $x$  intersects the real axis. We can “roll” this tangent line around the convex hull of the graph of  $\Lambda$ , to get all  $\Lambda^*$  values.
- ▶ Is the Legendre dual always convex?
- ▶ What is the Legendre dual of  $x^2$ ? Of the function equal to 0 at 0 and  $\infty$  everywhere else?
- ▶ How are derivatives of  $\Lambda$  and  $\Lambda^*$  related?

# Legendre transform

- ▶ Define **Legendre transform** (or Legendre dual) of a function  $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

- ▶ Let's describe the Legendre dual geometrically if  $d = 1$ :  $\Lambda^*(x)$  is where tangent line to  $\Lambda$  of slope  $x$  intersects the real axis. We can “roll” this tangent line around the convex hull of the graph of  $\Lambda$ , to get all  $\Lambda^*$  values.
- ▶ Is the Legendre dual always convex?
- ▶ What is the Legendre dual of  $x^2$ ? Of the function equal to 0 at 0 and  $\infty$  everywhere else?
- ▶ How are derivatives of  $\Lambda$  and  $\Lambda^*$  related?
- ▶ What is the Legendre dual of the Legendre dual of a convex function?

# Legendre transform

- ▶ Define **Legendre transform** (or Legendre dual) of a function  $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

- ▶ Let's describe the Legendre dual geometrically if  $d = 1$ :  $\Lambda^*(x)$  is where tangent line to  $\Lambda$  of slope  $x$  intersects the real axis. We can “roll” this tangent line around the convex hull of the graph of  $\Lambda$ , to get all  $\Lambda^*$  values.
- ▶ Is the Legendre dual always convex?
- ▶ What is the Legendre dual of  $x^2$ ? Of the function equal to 0 at 0 and  $\infty$  everywhere else?
- ▶ How are derivatives of  $\Lambda$  and  $\Lambda^*$  related?
- ▶ What is the Legendre dual of the Legendre dual of a convex function?
- ▶ What's the higher dimensional analog of rolling the tangent line?

Legendre transform

Large deviations

Legendre transform

Large deviations

## Recall: moment generating functions

- ▶ Let  $X$  be a random variable.

## Recall: moment generating functions

- ▶ Let  $X$  be a random variable.
- ▶ The **moment generating function** of  $X$  is defined by  $M(t) = M_X(t) := E[e^{tX}]$ .

## Recall: moment generating functions

- ▶ Let  $X$  be a random variable.
- ▶ The **moment generating function** of  $X$  is defined by  $M(t) = M_X(t) := E[e^{tX}]$ .

## Recall: moment generating functions

- ▶ Let  $X$  be a random variable.
- ▶ The **moment generating function** of  $X$  is defined by  $M(t) = M_X(t) := E[e^{tX}]$ .
- ▶ When  $X$  is discrete, can write  $M(t) = \sum_x e^{tx} p_X(x)$ . So  $M(t)$  is a weighted average of countably many exponential functions.

## Recall: moment generating functions

- ▶ Let  $X$  be a random variable.
- ▶ The **moment generating function** of  $X$  is defined by  $M(t) = M_X(t) := E[e^{tX}]$ .
- ▶ When  $X$  is discrete, can write  $M(t) = \sum_x e^{tx} p_X(x)$ . So  $M(t)$  is a weighted average of countably many exponential functions.
- ▶ When  $X$  is continuous, can write  $M(t) = \int_{-\infty}^{\infty} e^{tx} f(x) dx$ . So  $M(t)$  is a weighted average of a continuum of exponential functions.

## Recall: moment generating functions

- ▶ Let  $X$  be a random variable.
- ▶ The **moment generating function** of  $X$  is defined by  $M(t) = M_X(t) := E[e^{tX}]$ .
- ▶ When  $X$  is discrete, can write  $M(t) = \sum_x e^{tx} p_X(x)$ . So  $M(t)$  is a weighted average of countably many exponential functions.
- ▶ When  $X$  is continuous, can write  $M(t) = \int_{-\infty}^{\infty} e^{tx} f(x) dx$ . So  $M(t)$  is a weighted average of a continuum of exponential functions.
- ▶ We always have  $M(0) = 1$ .

## Recall: moment generating functions

- ▶ Let  $X$  be a random variable.
- ▶ The **moment generating function** of  $X$  is defined by  $M(t) = M_X(t) := E[e^{tX}]$ .
- ▶ When  $X$  is discrete, can write  $M(t) = \sum_x e^{tx} p_X(x)$ . So  $M(t)$  is a weighted average of countably many exponential functions.
- ▶ When  $X$  is continuous, can write  $M(t) = \int_{-\infty}^{\infty} e^{tx} f(x) dx$ . So  $M(t)$  is a weighted average of a continuum of exponential functions.
- ▶ We always have  $M(0) = 1$ .
- ▶ If  $b > 0$  and  $t > 0$  then  $E[e^{tX}] \geq E[e^{t \min\{X, b\}}] \geq P\{X \geq b\} e^{tb}$ .

## Recall: moment generating functions

- ▶ Let  $X$  be a random variable.
- ▶ The **moment generating function** of  $X$  is defined by  $M(t) = M_X(t) := E[e^{tX}]$ .
- ▶ When  $X$  is discrete, can write  $M(t) = \sum_x e^{tx} p_X(x)$ . So  $M(t)$  is a weighted average of countably many exponential functions.
- ▶ When  $X$  is continuous, can write  $M(t) = \int_{-\infty}^{\infty} e^{tx} f(x) dx$ . So  $M(t)$  is a weighted average of a continuum of exponential functions.
- ▶ We always have  $M(0) = 1$ .
- ▶ If  $b > 0$  and  $t > 0$  then  $E[e^{tX}] \geq E[e^{t \min\{X, b\}}] \geq P\{X \geq b\} e^{tb}$ .
- ▶ If  $X$  takes both positive and negative values with positive probability then  $M(t)$  grows at least exponentially fast in  $|t|$  as  $|t| \rightarrow \infty$ .

## Recall: moment generating functions for i.i.d. sums

- ▶ We showed that if  $Z = X + Y$  and  $X$  and  $Y$  are independent, then  $M_Z(t) = M_X(t)M_Y(t)$

## Recall: moment generating functions for i.i.d. sums

- ▶ We showed that if  $Z = X + Y$  and  $X$  and  $Y$  are independent, then  $M_Z(t) = M_X(t)M_Y(t)$
- ▶ If  $X_1 \dots X_n$  are i.i.d. copies of  $X$  and  $Z = X_1 + \dots + X_n$  then what is  $M_Z$ ?

## Recall: moment generating functions for i.i.d. sums

- ▶ We showed that if  $Z = X + Y$  and  $X$  and  $Y$  are independent, then  $M_Z(t) = M_X(t)M_Y(t)$
- ▶ If  $X_1 \dots X_n$  are i.i.d. copies of  $X$  and  $Z = X_1 + \dots + X_n$  then what is  $M_Z$ ?
- ▶ Answer:  $M_X^n$ .

# Large deviations

- ▶ Consider i.i.d. random variables  $X_i$ . Can we show that  $P(S_n \geq na) \rightarrow 0$  exponentially fast when  $a > E[X_i]$ ?

# Large deviations

- ▶ Consider i.i.d. random variables  $X_i$ . Can we show that  $P(S_n \geq na) \rightarrow 0$  exponentially fast when  $a > E[X_i]$ ?
- ▶ Kind of a quantitative form of the weak law of large numbers. The empirical average  $A_n$  is *very* unlikely to be  $\epsilon$  away from its expected value (where “very” means with probability less than some exponentially decaying function of  $n$ ).

# General large deviation principle

- ▶ More general framework: a *large deviation principle* describes limiting behavior as  $n \rightarrow \infty$  of family  $\{\mu_n\}$  of measures on measure space  $(\mathcal{X}, \mathcal{B})$  in terms of a *rate function*  $I$ .

# General large deviation principle

- ▶ More general framework: a *large deviation principle* describes limiting behavior as  $n \rightarrow \infty$  of family  $\{\mu_n\}$  of measures on measure space  $(\mathcal{X}, \mathcal{B})$  in terms of a *rate function*  $I$ .
- ▶ The **rate function** is a lower-semicontinuous map  $I : \mathcal{X} \rightarrow [0, \infty]$ . (The sets  $\{x : I(x) \leq a\}$  are closed — rate function called “good” if these sets are compact.)

# General large deviation principle

- ▶ More general framework: a *large deviation principle* describes limiting behavior as  $n \rightarrow \infty$  of family  $\{\mu_n\}$  of measures on measure space  $(\mathcal{X}, \mathcal{B})$  in terms of a *rate function*  $I$ .
- ▶ The **rate function** is a lower-semicontinuous map  $I : \mathcal{X} \rightarrow [0, \infty]$ . (The sets  $\{x : I(x) \leq a\}$  are closed — rate function called “good” if these sets are compact.)
- ▶ **DEFINITION:**  $\{\mu_n\}$  satisfy LDP with rate function  $I$  and speed  $n$  if for all  $\Gamma \in \mathcal{B}$ ,

$$- \inf_{x \in \Gamma^0} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq - \inf_{x \in \bar{\Gamma}} I(x).$$

# General large deviation principle

- ▶ More general framework: a *large deviation principle* describes limiting behavior as  $n \rightarrow \infty$  of family  $\{\mu_n\}$  of measures on measure space  $(\mathcal{X}, \mathcal{B})$  in terms of a *rate function*  $I$ .
- ▶ The **rate function** is a lower-semicontinuous map  $I : \mathcal{X} \rightarrow [0, \infty]$ . (The sets  $\{x : I(x) \leq a\}$  are closed — rate function called “good” if these sets are compact.)
- ▶ **DEFINITION:**  $\{\mu_n\}$  satisfy LDP with rate function  $I$  and speed  $n$  if for all  $\Gamma \in \mathcal{B}$ ,

$$-\inf_{x \in \Gamma^0} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq -\inf_{x \in \bar{\Gamma}} I(x).$$

- ▶ **INTUITION:** when “near  $x$ ” the probability density function for  $\mu_n$  is tending to zero like  $e^{-I(x)n}$ , as  $n \rightarrow \infty$ .

# General large deviation principle

- ▶ More general framework: a *large deviation principle* describes limiting behavior as  $n \rightarrow \infty$  of family  $\{\mu_n\}$  of measures on measure space  $(\mathcal{X}, \mathcal{B})$  in terms of a *rate function*  $I$ .
- ▶ The **rate function** is a lower-semicontinuous map  $I : \mathcal{X} \rightarrow [0, \infty]$ . (The sets  $\{x : I(x) \leq a\}$  are closed — rate function called “good” if these sets are compact.)
- ▶ **DEFINITION:**  $\{\mu_n\}$  satisfy LDP with rate function  $I$  and speed  $n$  if for all  $\Gamma \in \mathcal{B}$ ,

$$-\inf_{x \in \Gamma^0} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq -\inf_{x \in \bar{\Gamma}} I(x).$$

- ▶ **INTUITION:** when “near  $x$ ” the probability density function for  $\mu_n$  is tending to zero like  $e^{-I(x)n}$ , as  $n \rightarrow \infty$ .
- ▶ **Simple case:**  $I$  is continuous,  $\Gamma$  is closure of its interior.

# General large deviation principle

- ▶ More general framework: a *large deviation principle* describes limiting behavior as  $n \rightarrow \infty$  of family  $\{\mu_n\}$  of measures on measure space  $(\mathcal{X}, \mathcal{B})$  in terms of a *rate function*  $I$ .
- ▶ The **rate function** is a lower-semicontinuous map  $I : \mathcal{X} \rightarrow [0, \infty]$ . (The sets  $\{x : I(x) \leq a\}$  are closed — rate function called “good” if these sets are compact.)
- ▶ **DEFINITION:**  $\{\mu_n\}$  satisfy LDP with rate function  $I$  and speed  $n$  if for all  $\Gamma \in \mathcal{B}$ ,

$$- \inf_{x \in \Gamma^0} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq - \inf_{x \in \bar{\Gamma}} I(x).$$

- ▶ **INTUITION:** when “near  $x$ ” the probability density function for  $\mu_n$  is tending to zero like  $e^{-I(x)n}$ , as  $n \rightarrow \infty$ .
- ▶ **Simple case:**  $I$  is continuous,  $\Gamma$  is closure of its interior.
- ▶ **Question:** How would  $I$  change if we replaced the measures  $\mu_n$  by weighted measures  $e^{(\lambda n, \cdot)} \mu_n$ ?

# General large deviation principle

- ▶ More general framework: a *large deviation principle* describes limiting behavior as  $n \rightarrow \infty$  of family  $\{\mu_n\}$  of measures on measure space  $(\mathcal{X}, \mathcal{B})$  in terms of a *rate function*  $I$ .
- ▶ The **rate function** is a lower-semicontinuous map  $I : \mathcal{X} \rightarrow [0, \infty]$ . (The sets  $\{x : I(x) \leq a\}$  are closed — rate function called “good” if these sets are compact.)
- ▶ **DEFINITION:**  $\{\mu_n\}$  satisfy LDP with rate function  $I$  and speed  $n$  if for all  $\Gamma \in \mathcal{B}$ ,

$$- \inf_{x \in \Gamma^0} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq - \inf_{x \in \bar{\Gamma}} I(x).$$

- ▶ **INTUITION:** when “near  $x$ ” the probability density function for  $\mu_n$  is tending to zero like  $e^{-I(x)n}$ , as  $n \rightarrow \infty$ .
- ▶ **Simple case:**  $I$  is continuous,  $\Gamma$  is closure of its interior.
- ▶ **Question:** How would  $I$  change if we replaced the measures  $\mu_n$  by weighted measures  $e^{(\lambda n, \cdot)} \mu_n$ ?
- ▶ Replace  $I(x)$  by  $I(x) - (\lambda, x)$ ? What is  $\inf_x I(x) - (\lambda, x)$ ?

# Cramer's theorem

- ▶ Let  $\mu_n$  be law of empirical mean  $A_n = \frac{1}{n} \sum_{j=1}^n X_j$  for i.i.d. vectors  $X_1, X_2, \dots, X_n$  in  $\mathbb{R}^d$  with same law as  $X$ .

# Cramer's theorem

- ▶ Let  $\mu_n$  be law of empirical mean  $A_n = \frac{1}{n} \sum_{j=1}^n X_j$  for i.i.d. vectors  $X_1, X_2, \dots, X_n$  in  $\mathbb{R}^d$  with same law as  $X$ .
- ▶ Define **log moment generating function** of  $X$  by

$$\Lambda(\lambda) = \Lambda_X(\lambda) = \log M_X(\lambda) = \log \mathbb{E}e^{(\lambda, X)},$$

where  $(\cdot, \cdot)$  is inner product on  $\mathbb{R}^d$ .

# Cramer's theorem

- ▶ Let  $\mu_n$  be law of empirical mean  $A_n = \frac{1}{n} \sum_{j=1}^n X_j$  for i.i.d. vectors  $X_1, X_2, \dots, X_n$  in  $\mathbb{R}^d$  with same law as  $X$ .
- ▶ Define **log moment generating function** of  $X$  by

$$\Lambda(\lambda) = \Lambda_X(\lambda) = \log M_X(\lambda) = \log \mathbb{E}e^{(\lambda, X)},$$

where  $(\cdot, \cdot)$  is inner product on  $\mathbb{R}^d$ .

- ▶ Define **Legendre transform** of  $\Lambda$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

# Cramer's theorem

- ▶ Let  $\mu_n$  be law of empirical mean  $A_n = \frac{1}{n} \sum_{j=1}^n X_j$  for i.i.d. vectors  $X_1, X_2, \dots, X_n$  in  $\mathbb{R}^d$  with same law as  $X$ .
- ▶ Define **log moment generating function** of  $X$  by

$$\Lambda(\lambda) = \Lambda_X(\lambda) = \log M_X(\lambda) = \log \mathbb{E}e^{(\lambda, X)},$$

where  $(\cdot, \cdot)$  is inner product on  $\mathbb{R}^d$ .

- ▶ Define **Legendre transform** of  $\Lambda$  by

$$\Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}.$$

- ▶ **CRAMER'S THEOREM:**  $\mu_n$  satisfy LDP with convex rate function  $\Lambda^*$ .

## Thinking about Cramer's theorem

- ▶ Let  $\mu_n$  be law of empirical mean  $A_n = \frac{1}{n} \sum_{j=1}^n X_j$ .

# Thinking about Cramer's theorem

- ▶ Let  $\mu_n$  be law of empirical mean  $A_n = \frac{1}{n} \sum_{j=1}^n X_j$ .
- ▶ **CRAMER'S THEOREM:**  $\mu_n$  satisfy LDP with convex rate function

$$I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\},$$

where  $\Lambda(\lambda) = \log M(\lambda) = \mathbb{E}e^{(\lambda, X_1)}$ .

# Thinking about Cramer's theorem

- ▶ Let  $\mu_n$  be law of empirical mean  $A_n = \frac{1}{n} \sum_{j=1}^n X_j$ .
- ▶ **CRAMER'S THEOREM:**  $\mu_n$  satisfy LDP with convex rate function

$$I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\},$$

where  $\Lambda(\lambda) = \log M(\lambda) = \mathbb{E}e^{(\lambda, X_1)}$ .

- ▶ This means that for all  $\Gamma \in \mathcal{B}$  we have this **asymptotic lower bound** on probabilities  $\mu_n(\Gamma)$

$$- \inf_{x \in \Gamma^0} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma),$$

so (up to sub-exponential error)  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .

# Thinking about Cramer's theorem

- ▶ Let  $\mu_n$  be law of empirical mean  $A_n = \frac{1}{n} \sum_{j=1}^n X_j$ .
- ▶ **CRAMER'S THEOREM:**  $\mu_n$  satisfy LDP with convex rate function

$$I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\},$$

where  $\Lambda(\lambda) = \log M(\lambda) = \mathbb{E}e^{(\lambda, X_1)}$ .

- ▶ This means that for all  $\Gamma \in \mathcal{B}$  we have this **asymptotic lower bound** on probabilities  $\mu_n(\Gamma)$

$$-\inf_{x \in \Gamma^0} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma),$$

so (up to sub-exponential error)  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .

- ▶ and this **asymptotic upper bound** on the probabilities  $\mu_n(\Gamma)$

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq -\inf_{x \in \bar{\Gamma}} I(x),$$

which says (up to subexponential error)  $\mu_n(\Gamma) \leq e^{-n \inf_{x \in \bar{\Gamma}} I(x)}$ .

## Proving Cramer upper bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .

## Proving Cramer upper bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .
- ▶ For simplicity, assume that  $\Lambda$  is defined for all  $x$  (which implies that  $X$  has moments of all orders and  $\Lambda$  and  $\Lambda^*$  are strictly convex, and the derivatives of  $\Lambda$  and  $\Lambda'$  are inverses of each other). It is also enough to consider the case  $X$  has mean zero, which implies that  $\Lambda(0) = 0$  is a minimum of  $\Lambda$ , and  $\Lambda^*(0) = 0$  is a minimum of  $\Lambda^*$ .

## Proving Cramer upper bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .
- ▶ For simplicity, assume that  $\Lambda$  is defined for all  $x$  (which implies that  $X$  has moments of all orders and  $\Lambda$  and  $\Lambda^*$  are strictly convex, and the derivatives of  $\Lambda$  and  $\Lambda'$  are inverses of each other). It is also enough to consider the case  $X$  has mean zero, which implies that  $\Lambda(0) = 0$  is a minimum of  $\Lambda$ , and  $\Lambda^*(0) = 0$  is a minimum of  $\Lambda^*$ .
- ▶ We aim to show (up to subexponential error) that  $\mu_n(\Gamma) \leq e^{-n \inf_{x \in \Gamma} I(x)}$ .

## Proving Cramer upper bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .
- ▶ For simplicity, assume that  $\Lambda$  is defined for all  $x$  (which implies that  $X$  has moments of all orders and  $\Lambda$  and  $\Lambda^*$  are strictly convex, and the derivatives of  $\Lambda$  and  $\Lambda^*$  are inverses of each other). It is also enough to consider the case  $X$  has mean zero, which implies that  $\Lambda(0) = 0$  is a minimum of  $\Lambda$ , and  $\Lambda^*(0) = 0$  is a minimum of  $\Lambda^*$ .
- ▶ We aim to show (up to subexponential error) that  $\mu_n(\Gamma) \leq e^{-n \inf_{x \in \Gamma} I(x)}$ .
- ▶ If  $\Gamma$  were singleton set  $\{x\}$  we could find the  $\lambda$  corresponding to  $x$ , so  $\Lambda^*(x) = (\lambda, x) - \Lambda(\lambda)$ . Note then that

$$\mathbb{E}e^{(n\lambda, A_n)} = \mathbb{E}e^{(\lambda, S_n)} = M_X^n(\lambda) = e^{n\Lambda(\lambda)},$$

and also  $\mathbb{E}e^{(n\lambda, A_n)} \geq e^{n(\lambda, x)} \mu_n\{x\}$ . Taking logs and dividing by  $n$  gives  $\Lambda(\lambda) \geq \frac{1}{n} \log \mu_n + (\lambda, x)$ , so that  $\frac{1}{n} \log \mu_n(\Gamma) \leq -\Lambda^*(x)$ , as desired.

## Proving Cramer upper bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .
- ▶ For simplicity, assume that  $\Lambda$  is defined for all  $x$  (which implies that  $X$  has moments of all orders and  $\Lambda$  and  $\Lambda^*$  are strictly convex, and the derivatives of  $\Lambda$  and  $\Lambda'$  are inverses of each other). It is also enough to consider the case  $X$  has mean zero, which implies that  $\Lambda(0) = 0$  is a minimum of  $\Lambda$ , and  $\Lambda^*(0) = 0$  is a minimum of  $\Lambda^*$ .
- ▶ We aim to show (up to subexponential error) that  $\mu_n(\Gamma) \leq e^{-n \inf_{x \in \Gamma} I(x)}$ .
- ▶ If  $\Gamma$  were singleton set  $\{x\}$  we could find the  $\lambda$  corresponding to  $x$ , so  $\Lambda^*(x) = (x, \lambda) - \Lambda(\lambda)$ . Note then that

$$\mathbb{E}e^{(n\lambda, A_n)} = \mathbb{E}e^{(\lambda, S_n)} = M_X^n(\lambda) = e^{n\Lambda(\lambda)},$$

and also  $\mathbb{E}e^{(n\lambda, A_n)} \geq e^{n(\lambda, x)} \mu_n\{x\}$ . Taking logs and dividing by  $n$  gives  $\Lambda(\lambda) \geq \frac{1}{n} \log \mu_n + (\lambda, x)$ , so that  $\frac{1}{n} \log \mu_n(\Gamma) \leq -\Lambda^*(x)$ , as desired.

- ▶ General  $\Gamma$ : cut into finitely many pieces, bound each piece?

## Proving Cramer lower bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .

## Proving Cramer lower bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .
- ▶ We aim to show that asymptotically  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .

## Proving Cramer lower bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .
- ▶ We aim to show that asymptotically  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .
- ▶ It's enough to show that for each given  $x \in \Gamma^0$ , we have that asymptotically  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .

## Proving Cramer lower bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .
- ▶ We aim to show that asymptotically  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .
- ▶ It's enough to show that for each given  $x \in \Gamma^0$ , we have that asymptotically  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .
- ▶ Idea is to weight the law of  $X$  by  $e^{(\lambda, x)}$  for some  $\lambda$  and normalize to get a new measure whose expectation is this point  $x$ . In this new measure,  $A_n$  is “typically” in  $\Gamma$  for large  $\Gamma$ , so the probability is of order 1.

## Proving Cramer lower bound

- ▶ Recall that  $I(x) = \Lambda^*(x) = \sup_{\lambda \in \mathbb{R}^d} \{(\lambda, x) - \Lambda(\lambda)\}$ .
- ▶ We aim to show that asymptotically  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .
- ▶ It's enough to show that for each given  $x \in \Gamma^0$ , we have that asymptotically  $\mu_n(\Gamma) \geq e^{-n \inf_{x \in \Gamma^0} I(x)}$ .
- ▶ Idea is to weight the law of  $X$  by  $e^{(\lambda, x)}$  for some  $\lambda$  and normalize to get a new measure whose expectation is this point  $x$ . In this new measure,  $A_n$  is “typically” in  $\Gamma$  for large  $\Gamma$ , so the probability is of order 1.
- ▶ But by how much did we have to modify the measure to make this typical? Not more than by factor  $e^{-n \inf_{x \in \Gamma^0} I(x)}$ .