# Information Theoretic Lower Bounds

**Recap**

We studied convex (stochastic) optimization problems

$$\text{minimize}_{\theta \in \Theta} \quad R(\theta), \qquad \Theta \subseteq \mathbb{R}^d \text{ convex, compact,} \qquad R: \Theta \to \mathbb{R} \text{ convex.}$$

**Def** (subgradients) The subgradient set of $R$ at $\theta$ is $\partial R(\theta) := \{g : R(\theta') \geq R(\theta) + \langle g, \theta'-\theta \rangle \ \forall \theta'\}$

cf. $\partial R(\theta)$ is convex and compact. subgradients are supporting hyperplanes of the epigraph



FOC: $\theta^* = \text{argmin}_{\theta \in \Theta} R(\theta)$   iff   $\langle g, \theta-\theta^* \rangle \geq 0 \quad \exists g \in \partial R(\theta^*) \quad \forall \theta \in \Theta$

**Example**
$$h(\theta) = |\theta - c| \qquad \partial h(\theta) = \begin{cases} \text{sign}(\theta - c) & \text{if } \theta \neq c \\ [-1, 1] & \text{if } \theta = c \end{cases}$$

cf. Useful for e.g. $\min \frac{1}{n}\sum \ell(\theta; X, Y) + \lambda \|\theta\|_1$   or   $\min \frac{1}{n}\sum |Y_i - \theta^T X_i|$

Stochastic (sub)gradient descent : Consider stochastic subgradients $g(\theta; \xi)$ s.t. for some independent RV $\xi$, $\mathbb{E} \, g(\theta; \xi) \in \partial R(\theta)$.

$$\theta^{k+1} \leftarrow \Pi_\Theta \left( \theta^k - \alpha_k \, g(\theta^k; \xi_k) \right) \qquad : \text{SGD update}$$

e.g. $\xi_k = (X_k, Y_k). \qquad g(\theta; \xi) = \nabla_\theta \ell(\theta, X, Y).$   (or you can use batch $> 1$)

**Theorem** Let $\theta^* \in \text{argmin}_{\theta \in \Theta} R(\theta) > -\infty$, $D > 0$ s.t. $\sup_{\theta \in \Theta} \|\theta - \theta^*\|_2 \leq D$, $\exists M > 0$ s.t. $\sup_\theta \mathbb{E} \|g(\theta; \xi)\|_2^2 \leq M^2$.
Let $\alpha_k$ be dec, pos. step sizes, and $\bar{\theta}_K = \frac{1}{K} \sum^K \theta_k$

$$\mathbb{E}[R(\bar{\theta}_K) - R(\theta^*)] \leq \frac{D^2}{2K\alpha_K} + \frac{1}{2K} \sum_1^k \alpha_k M^2$$

Pf) Identical as last week's result for differentiable functions.

**Cor** Setting $\alpha_k = \frac{D}{M\sqrt{k}}$, $\mathbb{E} R(\bar{\theta}_K) - R(\theta^*) \leq \frac{3DM}{2\sqrt{K}}$.

**Rmk** To solve problem to $\varepsilon$-accuracy, you need $O(\frac{1}{\varepsilon^2})$ - iterations of SGD. Compared to $\log(1/\varepsilon)$ of IPMs, this is terribly slow.

**Q** Can we improve the SGD rate of conv? Or is it "optimal"? i.e. unimprovable

# Minimax rates
Considers worst-case optimality gap among a class of "problems".

<u>Components</u>
1) A collection of functions $\mathcal{F} := \{\theta \mapsto \mathbb{E}_P \ell(\theta; z) : P \in \mathcal{P} \ \forall\}$ induced by class of probabilities $\mathcal{P}$
                                                                      indep. of everything
2) Closed, convex set $\Theta \subseteq \mathbb{R}^d$

3) A stochastic gradient oracle $g : \mathbb{R}^d \times \Xi \times \mathcal{F} \to \mathbb{R}^d$   s.t.   $\mathbb{E} g(\theta; z; P) \in \partial R(\theta)$. Implicitly a prob. distribution on $\Xi$ induced by $P$.

    ↳ Think of this as # access to data points.

<u>Example</u>
$R_P(\theta) = \mathbb{E}_P \log(1 + \exp(-Y\theta^T X))$   : Logistic regression     $\mathcal{P} = \{P \text{ s.t. } Y \in \{-1, 1\}, X \in [0,1]^d \text{ a.s.}\}$
$\mathcal{F} = \{R_P : P \in \mathcal{P}\}$      Define $z = (X, Y)$.
$g(\theta; z; P) = \nabla_\theta \ell(\theta; X, Y) = -\dfrac{YX}{1 + \exp(Y\theta^T X)}$,    $\mathbb{E}_{z \sim P} g(\theta; z; P) = \nabla_\theta R_P(\theta)$.

View algorithms as making $K$ queries to the stoch. grad. oracle at $\theta_1, \ldots, \theta_K$ (adaptively) and it returns $\hat{\theta}_K$.

Optimality gap: $\mathbb{E}_{z_i}\left[ R_P(\hat{\theta}(z_1, \ldots, z_K)) - \inf_{\theta \in \Theta} R_P(\theta) \right]$        ↙ IDK what this is. So we think of nature choosing the worst-possible instance.

Measure performance w.r.t. hardest problem (dist. $P \in \mathcal{P}$)
$$\sup_{\ell, P \in \mathcal{P}} \mathbb{E}_{z_i}\left[ R_P(\hat{\theta}(z_1, \ldots, z_K)) - \inf_{\theta \in \Theta} R_P(\theta) \right] \qquad \cdots (*)$$
i.e. I want $\hat{\theta}$ to achieve low error uniformly over $P \in \mathcal{P}$.

My job as statistical modeler / stoch optimizer is to come up with algo $\hat{\theta}$ st. $(*)$ low.

The optimal algo gives the <span style="color:red">minimax risk</span>
$$\mathcal{M}_K(\Theta, \mathcal{P}) := \inf_{\hat{\theta}} \sup_{\ell, P \in \mathcal{P}} \mathbb{E}_{z_i}\left[ R_P(\hat{\theta}(z_1, \ldots, z_K)) - \inf_{\theta \in \Theta} R_P(\theta) \right]$$

where inf is over all measurable ftns of $z_1, \ldots, z_K$.

<u>Rmk</u> We can consider random procedures that incorporate more randomness using same techniques / bounds / proofs.

This measures extent to which we can optimize worst-case opt gap over the class of problems $P \in \mathcal{P}$, using finite access to stoch grad oracle.

<u>Rmk</u> Similar quantities can be defined for estimation problems where $z_1, \ldots, z_K$ are data points. This would measure sample complexity. Everything in today's class has a straightforward analogue to this setting.

<u>Rmk</u> Though we define minimax risk over $K$ queries to stoch grad oracle, we can see same bound applies to all procedures using $K$ samples.

<u>Rmk</u>   Minimax can be very conservative.

<u>Step 0</u>  (Worst-case → Bayesian)   Let $\{P_v\} \subseteq P$ be a collection of distributions indexed by finite or countable $\mathcal{V}$. Let $\pi$ be a prob on $\mathcal{V}$. For any fixed $\ell$,

$$\sup_{P \in \mathcal{P}} \mathbb{E}\left[R_P(\hat{\theta}_K) - \inf_{\theta \in \Theta} R_P(\theta)\right] \geq \sum_{v \in V} \pi(v) \underbrace{\mathbb{E}\left[R_{P_v}(\hat{\theta}_K) - \inf_{\theta \in \Theta} R_{P_v}(\theta)\right]}_{R_v}$$
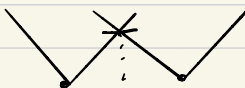
<u>Step 1</u>  (Reduction from optimization to hypothesis testing)

<u>Def</u>   For two cvx functions $R_0, R_1$ define the optimization separation bet $R_0$ & $R_1$ as

$$d_{opt}(R_0, R_1) := \sup\left\{\delta \geq 0 : \begin{array}{l} R_1(\theta) \leq R_1^* + \delta \Rightarrow R_0(\theta) \geq R_0^* + \delta \\ R_0(\theta) \leq R_0^* + \delta \Rightarrow R_1(\theta) \geq R_1^* + \delta \end{array} \text{ for any } \theta \in \Theta \right\}$$

where $R_v^* = \inf_{\theta \in \Theta} R_v(\theta)$.

    ↳ <span style="color:blue">Optimizing $R_1$ means we can't optimize $R_0$ very well  if $d_{opt}(R_0, R_1)$ large. This says if we optimized one function well, then we couldn't have optimized other functions well-separated in $d_{opt}$.</span>

<u>e.g.</u>  If $\theta$ is s.t. $R_1(\theta) - R_1^* \leq d_{opt}(R_0, R_1)$, then $\theta$ cannot optimize $R_0$ well.

<u>Ex</u>   $R_0(\theta) = |\theta - \nabla|$   $R_1(\theta) = |\theta + \nabla|$                             $d_{opt}(R_0, R_1) = \nabla$.

Consider canonical hypothesis testing problem:

   1) Nature chooses $V \in \mathcal{V}$ unif. at random

   2) Cond. on $V = v$, we observe subgradients associated with $R_{P_v}$ for i.i.d. $\xi_1, \dots, \xi_K$.

<u>Goal</u>   Figure out which index Nature chose.

    If we can opt to better than $d_{opt}(R_v, R_{v'})$ $\forall v \in V$, then we can identify $V = v$.

<u>Lemma</u>   Let $V$ be drawn u.a.r. from $\mathcal{V}$, $|V| < \infty$, and assume $\{P_v\}_{v \in V}$ is $\delta$-separated:
$d_{opt}(R_v, R_{v'}) \geq \delta$  $\forall v \neq v' \in \mathcal{V}$. For any fixed $\ell$, <span style="color:blue;">over $V$ & $\xi_1^K$</span>

Then, $\quad \dfrac{1}{|V|} \sum_{v \in V} \mathbb{E}\left[R_v(\hat{\theta}_K) - R_v^*\right] \geq \delta \inf_{\hat{v}} \mathbb{P}(\hat{v} \neq V)$

where inf is over all testing procedures based on observed data.
  <span style="color:blue;">↳ In particular, this lower bounds Bayes-risk with $\pi(v) = \frac{1}{|V|}$.</span>

<u>Game plan</u>   Construct a class of well-separated functions in $d_{opt}$. Then show testing among them is difficult.
  ↳ Trade-off: Easier to distinguish functions when $\delta$ large, and vice-versa.

Pf) $\mathbb{E}[R_v(\hat{\theta}) - R_i^*] \geq \delta \, \mathbb{E} \, \mathbb{1}\{R_v(\hat{\theta}) - R_i^* \geq \delta\} = \delta \, \mathbb{P}_v(R_v(\hat{\theta}) - R_i^* \geq \delta)$. $\checkmark$ $\mathbb{P}_v$ is over $\overline{\mathfrak{z}}^K$, for the oracle it defines.

Define the hypothesis test $\hat{v} = \begin{cases} v & \text{if } R_v(\hat{\theta}) \leq R_i^* + \delta \\ \text{random} & \text{o/w} \leftarrow \text{this can be arbitrary} \end{cases}$

$\llcorner$ This is well-defined since $d_{opt}(R_v, R_{v'}) \geq \delta \quad \forall v \neq v'$, so such $v$ is unique.

Now note $\hat{v} \neq v \implies R_v(\hat{\theta}) \geq R_i^* + \delta$. So $\mathbb{P}(\hat{v} \neq v) \leq \mathbb{P}_v(R_v(\hat{\theta}) \geq R_i^* + \delta)$.

Hence, $\frac{1}{|V|} \sum_{v \in V} \mathbb{E}[R_v(\hat{\theta}_K) - R_i^*] \geq \delta \frac{1}{|V|} \sum_{v \in V} \mathbb{P}_v(\hat{v} \neq v) = \delta \cdot \mathbb{P}(\hat{v} \neq V)$ by def of $V$. $\boxtimes$

# Le Cam's method.   Reduction to binary hypothesis testing.   $V = \{-1, +1\}$

We construct $P_1, P_{-1}$, and show hardness of optimizing on $\mathbb{R}^1 \supseteq \Theta$.

Def
| TV distance | $\|P - Q\|_{TV} = \sup\limits_{A \subseteq \Xi \text{ meas.}} |P(A) - Q(A)|$ |
| KL divergence | $D_{KL}(P, Q) = \int p(\mathfrak{z}) \log \frac{p(\mathfrak{z})}{q(\mathfrak{z})} \, d\mu(\mathfrak{z})$ |

Lemma 1  $1 - \|P_1 - P_{-1}\|_{TV} = \inf\limits_{\hat{v}} \{P_1(\hat{v} \neq 1) + P_{-1}(\hat{v} \neq -1)\}$ $\leftarrow$ best prob of being wrong in binary hypo test.

Pf) Any $\hat{v}: \Xi \to \{-1, 1\}$, define $A = \hat{v}^{-1}\{1\}$, $A^c = \hat{v}^{-1}\{-1\}$ so $P_1(\hat{v} \neq 1) + P_{-1}(\hat{v} \neq -1) = P_1(A^c) + P_{-1}(A) = 1 - P_1(A) + P_{-1}(A)$

Taking inf over $\hat{v}$, RHS $= \inf\limits_{A \subseteq \Xi \text{ meas}} \{1 - P_1(A) + P_{-1}(A)\} = 1 - \sup\limits_{A \subseteq \Xi \text{ meas}} P_1(A) - P_{-1}(A) = 1 - \|P_1 - P_{-1}\|_{TV}$. $\boxtimes$

We get $M_n(\Theta, \mathcal{P}) \geq \inf\limits_{\hat{\theta}_K} \max\limits_{v \in \{\pm 1\}} \mathbb{E}_{P_v}[R_v(\hat{\theta}_K) - R_i^*] \geq \delta \cdot \frac{1}{2} \inf\limits_{\hat{v}} \{P_1^K(\hat{v} \neq 1) + P_{-1}^K(\hat{v} \neq -1)\}$

where $P_{\pm 1}^K$ are $K$-product distr over $\mathfrak{z}_1, ..., \mathfrak{z}_K$.   $= \delta \cdot \frac{1}{2} \cdot (1 - \|P_1^K - P_{-1}^K\|_{TV})$

$\llcorner$ If $R_1, R_{-1}$ are diff, then $\delta$ is big (optimizing one makes other worse), but this will likely make $P_1, P_{-1}$ to be different.

Now, $\|P_1^K - P_{-1}^K\|_{TV}$ is unwieldy since TV distance don't play well with products.

Lemma 2  (Pinsker)    $\|P - Q\|_{TV}^2 \leq \frac{1}{2} D_{KL}(P, Q)$

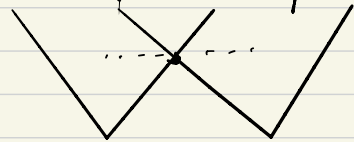Lemma 3  $D_{KL}(P_1^K, P_{-1}^K) = K \cdot D_{KL}(P_1, P_{-1})$

Pf) $\int \prod_{h=1}^K P_1(\mathfrak{z}_h) \log \frac{\prod_{h=1}^K P_1(\mathfrak{z}_h)}{\prod_{h=1}^K P_{-1}(\mathfrak{z}_h)} \, d\mathfrak{z}_1 \cdots d\mathfrak{z}_K = \sum_{h=1}^K \int \left(\prod_{h=1}^K P_h(\mathfrak{z}_h)\right) \log \frac{P_1(\mathfrak{z}_h)}{P_{-1}(\mathfrak{z}_h)} \, d\mathfrak{z}_1 \cdots d\mathfrak{z}_K$

$= \sum_{h=1}^K \int P_h(\mathfrak{z}_h) \log \frac{P_1(\mathfrak{z}_h)}{P_{-1}(\mathfrak{z}_h)} \, d\mathfrak{z}_h = \sum_{h=1}^K D_{KL}(P_1, P_{-1})$  $\boxtimes$.

So we get $\|P_1^K - P_{-1}^K\|_{TV} \leq \sqrt{\frac{K}{2} D_{KL}(P_1, P_{-1})}$. Plugging this in,

$M_K(\Theta, \mathcal{P}) \geq \frac{1}{2} d_{opt}(R_1, R_{-1}) \cdot \left(1 - \sqrt{\frac{K}{2} D_{KL}(P_1, P_{-1})}\right)$.  $\boxtimes$

Now, construct $P_1, P_{-1}$ (with corresponding oracles) s.t. they're close in KL so testing is hard, but well-separated in $d_{opt}$.

Let $\Theta \subseteq \mathbb{R}^d$ be s.t. $L^2$-ball of radius $D$ is included in $\Theta$. We consider stoch. grad. oracles s.t. $\mathbb{E}\| g(\theta;\xi;R_v)\|_2^2 \leq M^2$ for all $\theta \in \Theta$. Let $P$ be s.t. $\theta \mapsto \mathbb{E}_P(\theta)$ is $M$-Lip in $\|\cdot\|_2$.



**Construction**

Consider $d=1$ w.l.o.g.: Assume $[-D,D] \subseteq \Theta \subseteq \mathbb{R}$.

Define $R_1, R_{-1}$ s.t. $R_1(\theta) = \delta M|\theta - D|$   $R_{-1}(\theta) = \delta M|\theta + D|$.

From picture, $d_{opt}(R_1, R_{-1}) = \delta M D$.

Now, construct stoch. grad. oracle s.t. $\mathbb{E}\, g(\theta;\xi;R_v)^2 \leq M^2$.

For $\delta < 1$, the oracle for $R_v$, $v \in \{\pm 1\}$ is

flip coin $\xi \sim Ber\left(\frac{1+v\delta}{2}\right)$,   $g(\theta;\xi;R_v) = \begin{cases} vM \, sgn(\theta - vD) & \text{if } \xi = 1 \\ -vM \, sgn(\theta - vD) & \text{if } \xi = 0 \end{cases}$

cf. If $\theta = vD$, return random $\pm M$.

$\mathbb{E}_{R_v} g(\theta;\xi;R_v) = \frac{1+v\delta}{2} vM \, sgn(\theta - vD) - \frac{1-v\delta}{2} vM \, sgn(\theta - vD) = v\delta \cdot vM \, sgn(\theta - vD) = \delta M \, sgn(\theta - vD) \in \partial R_v(\theta)$

**Distance**

$D_{KL}(P_1, P_{-1}) = \sum_{\xi \in \{0,1\}} P_1(\xi) \log \frac{P_1(\xi)}{P_{-1}(\xi)} = \frac{1+\delta}{2} \log \frac{\frac{1+\delta}{2}}{\frac{1-\delta}{2}} + \frac{1-\delta}{2} \log \frac{\frac{1-\delta}{2}}{\frac{1+\delta}{2}}$

$= \delta \log \frac{1+\delta}{1-\delta}$

Taylor expansion of $x \mapsto \log \frac{1+x}{1-x} = a(x)$ at $x=0$:   $a'(x) = \frac{1}{1+x} + \frac{1}{1-x}$   $a''(x) = -\frac{1}{(1+x)^2} + \frac{1}{(1-x)^2}$

$\log \frac{1+\delta}{1-\delta} = a(0) + a'(0)\delta + \frac{a''(\hat\delta)}{2}\delta^2$   for some $\hat\delta \in [0,\delta]$

$\qquad\quad \leq \quad 0 \quad + 2\delta + 2\delta^2 \qquad$ if $\delta \leq \frac{1}{2}$   $\left( \because a''(\hat\delta) \leq \frac{1}{(1-\hat\delta)^2} \leq \frac{1}{(1-\frac12)^2} = 4 \right)$

$\qquad\quad \leq 3\delta \qquad\quad$ if $\delta \leq \frac{1}{2}$.

So $D_{KL}(P_1, P_{-1}) \leq 3\delta^2$.   Plugging this into Le Cam's result,

$$\mathcal{M}_K(\Theta, P) \geq \frac{1}{2} \delta M D \cdot \left( 1 - \sqrt{\frac{K}{2} 3\delta^2} \right).$$

**Set $\delta$**

I need $1 - \sqrt{\frac{K\delta^2}{2} 3}$ to behave like a const. i.e. $\delta^2 \sim \frac{1}{K}$.

Set $\delta = \frac{1}{\sqrt{6K}}$ so that $1 - \sqrt{\frac{3K}{2}\delta^2} = \frac{1}{2}$   $(K \geq 2)$.

Conclude $\mathcal{M}_K(\Theta, P) \geq \frac{MD}{4\sqrt{6K}}$.   for $K \geq 2$.

Recalling prev result, if $D_{inner} \mathbb{B}_2 \subseteq \Theta \subseteq D_{outer} \mathbb{B}_2$,

$$\frac{D_{inner} M}{\sqrt{K}} \lesssim \mathcal{M}_K(\Theta, P) \lesssim \frac{D_{outer} M}{\sqrt{K}}.$$

**Q** My construction was in $\mathbb{R}^1$. Surely things are harder in $d>1$. Can we show a more explicit bound with dimension?

# Assouad's method

Reduce opt to $d$ binary hypothesis tests.

Let $V = \{\pm 1\}^d$ be the $d$-dim binary hypercube. For each $v \in V$, we construct function $R_v$ and $P_v$, a joint distribution over $\mathcal{Z}$.

W.l.o.g. assume $0 \in \Theta$. We say $\{P_v\}_{v \in V}$ is $\delta$-separated in the Hamming distance if for all $\theta \in \Theta$

$$R_v(\theta) - R_v^* \geq \sum_{j=1}^d \delta_j \, \mathbb{1}\{\mathrm{sgn}(\theta_j) \neq v_j\}.$$

**Ex** $R_v(\theta) = \sum_{j=1}^d \delta_j |\theta_j - v_j|$, $V \subseteq \Theta$. So $R_v^* = 0 \; \forall v \in V$. Then, $R_v(\theta) - R_v^* = \sum_{j=1}^d \delta_j |\theta_j - v_j| \geq \sum_{j=1}^d \delta_j \mathbb{1}\{\mathrm{sgn}(\theta_j) \neq v_j\}$

**Lemma** (Assouad) Let $\{P_v\}_{v \in V}$ be $\delta$-separated in Hamming distance, $V = \{\pm 1\}^d$, Let $\bar{P}_{\pm j} = \frac{1}{2^{d-1}} \sum_{v_j = \pm 1} P_v$

$$\frac{1}{2^n} \sum_{v \in \{\pm 1\}^d} \mathbb{E}\, R_v(\hat{\theta}_k) - R_v^* \geq \frac{1}{2} \sum_i^d \delta_j \, (1 - \|\bar{P}_{+j} - \bar{P}_{-j}\|_{TV})$$

**Pf)** $\frac{1}{2^d} \sum_{v \in \{\pm 1\}^d} \mathbb{E}\left[ R_v(\hat{\theta}_k) - R_v^* \right] \geq \frac{1}{2^d} \sum_{v \in \{\pm 1\}^d} \sum_{j=1}^d \delta_j \, P_v\left(\mathrm{sgn}(\hat{\theta}_j) \neq v_j\right) = \sum_{j=1}^d \delta_j \; \frac{1}{2^d} \sum_{v \in \{\pm 1\}^d} P_v\left(\mathrm{sgn}(\hat{\theta}_j) \neq v_j\right)$

$$= \sum_i^d \delta_j \; \frac{1}{2^d} \left\{ \sum_{v_j = 1} P_v\left(\mathrm{sgn}(\hat{\theta}_j) \neq 1\right) + \sum_{v_j = -1} P_v\left(\mathrm{sgn}(\hat{\theta}_j) \neq -1\right) \right\}$$

$$= \frac{1}{2} \sum_i^d \delta_j \left\{ \bar{P}_{+j}\left(\mathrm{sgn}(\hat{\theta}_j) \neq 1\right) + \bar{P}_{-j}\left(\mathrm{sgn}(\hat{\theta}_j) \neq -1\right) \right\}$$

$$\geq \frac{1}{2} \sum_i^d \delta_j \left(1 - \|\bar{P}_{+j} - \bar{P}_{-j}\|_{TV}\right) \qquad \boxtimes.$$

**Cor** Let $P_{v, \pm j}$ be $P_{v'}$ where $v'$ is $v$ with $j^{\text{th}}$ entry of $v$ forced to be $\pm 1$. Then,

$$\mathcal{M}_K(\Theta, \mathcal{P}) \geq \frac{d\delta}{2}\left(1 - \sqrt{\max_{v \in V, 1 \leq j \leq d} D_{KL}\left(P_{v,+j}, P_{v,-j}\right)/2}\right) \qquad \text{whenever } \delta \cdot \mathbb{1}\text{-separation in Hamming holds.}$$

**Pf)** Note that $\bar{P}_{\pm j} = \frac{1}{2^d} \sum_v P_{v, \pm j}$ since $\sum_v P_{v, +j} = \sum_{v: v_j = 1} P_v + \sum_{v: v_j = -1} P_{v, +j} = 2 \sum_{v: v_j = 1} P_v$

Then from triangle inequality,

$$\|\bar{P}_{+j} - \bar{P}_{-j}\|_{TV} = \left\| \frac{1}{2^d} \sum_v (P_{v,+j} - P_{v,-j}) \right\|_{TV} \leq \frac{1}{2^d} \sum_v \|P_{v,+j} - P_{v,-j}\|_{TV} \leq \max_{v,j} \|P_{v,+j} - P_{v,-j}\|_{TV} \leq \max_{v,j} \sqrt{\frac{1}{2} D_{KL}(P_{v,j}, P_{v,-j})} \; \text{by Pinsker.}$$

$$\mathcal{M}_K(\Theta, \mathcal{P}) \geq \frac{d\delta}{2}\left(1 - \frac{1}{d}\sum_i^d \|\bar{P}_{+j} - \bar{P}_{-j}\|_{TV}\right) \geq \frac{d\delta}{2}\left(1 - \max_{v,j} \sqrt{\frac{1}{2} D_{KL}(P_{v,j}, P_{v,-j})}\right). \qquad \boxtimes.$$

Consider $[-v, v]^d \subseteq \Theta$, and $\mathbb{E}\|g(\theta; \mathcal{Z}; P)\|_i^2 \leq M^2 \;\; \forall \theta \in \Theta$, and $|R_P(\theta) - R_P(\theta')| \leq M \|\theta - \theta'\|_\infty$

**Construction** Let $R_v(\theta) = \frac{M\delta}{d} \|\theta - \mathcal{D} \cdot v\|_1$ for $v \in \{\pm 1\}^d$. Then

$$R_v(\theta) - R_v^* = R_v(\theta) \geq \frac{M\delta}{d} \sum_i^d D \mathbb{1}\{\mathrm{sgn}(\theta_j) \neq v_j\}. \text{ So } \frac{DM\delta}{d} \mathbb{1}\text{-separated in the Hamming distance.}$$

For stoch. grad. oracle, take $\mathcal{Z}_v = \begin{cases} e_j & \text{w.p. } \frac{1 + v_j \delta}{2d} \\ -e_j & \text{w.p. } \frac{1 - v_j \delta}{2d} \end{cases} \quad \forall j = 1, \dots, d$

$g(\theta; \mathcal{Z}; R_v) = \begin{cases} M \cdot v_j \cdot \mathrm{sgn}(\theta_j - \mathcal{D} v_j) \cdot e_j & \text{if } \mathcal{Z}_v = e_j \\ -M v_j \, \mathrm{sgn}(\theta_j - \mathcal{D} v_j) e_j & \text{if } \mathcal{Z}_v = -e_j \end{cases}$, $\pm M$ randomly if $\theta_j = \mathcal{D} v_j$

$\mathbb{E}_{R_v} g(\theta; \mathcal{Z}; R_v) = \sum_{j=1}^d \left\{ \frac{1 + v_j \delta}{2d} M v_j \, \mathrm{sgn}(\theta_j - \mathcal{D} v_j) e_j - \frac{1 - v_j \delta}{2d} M v_j \, \mathrm{sgn}(\theta_j - \mathcal{D} v_j) e_j \right\}$

$= \sum_{j=1}^d \frac{\delta M}{d} v_j^2 \, \mathrm{sgn}(\theta_j - \mathcal{D} v_j) e_j = \frac{\delta M}{d} \mathrm{sgn}(\theta - \mathcal{D} v) \in \partial R_v(\theta)$

We need to bound $D_{KL}(P_v, P_{v'})$ for $v \& v'$ that differ only in a single coordinate. W.l.o.g. let this be the first coordinate, and recall $P_v$ is a joint distribution over $\mathfrak{z}_i^k$. Since $P_v$ is a product distribution over i.i.d. $\mathfrak{z}_i$, let $\mathfrak{z}_1 \sim P_{v,\mathfrak{z}_1}$. Then $P_v = P_{v,\mathfrak{z}_1}^k$.

So $D_{KL}(P_v, P_{v'}) = D_{KL}(P_{v,\mathfrak{z}_1}^k, P_{v',\mathfrak{z}_1}^k) = k \, D_{KL}(P_{v,\mathfrak{z}_1}, P_{v',\mathfrak{z}_1})$

Now, note from definition that if $v_j = v'_j \; \forall j \neq 1$, with $v_1 = 1$, $v'_1 = -1$

$$D_{KL}(P_{v,\mathfrak{z}_1}, P_{v',\mathfrak{z}_1}) = P(\mathfrak{z}_v = e_1) \cdot \log \frac{P(\mathfrak{z}_v = e_1)}{P(\mathfrak{z}_{v'} = e_1)} + P(\mathfrak{z}_v = -e_1) \log \frac{P(\mathfrak{z}_v = -e_1)}{P(\mathfrak{z}_{v'} = -e_1)}$$

$$= \frac{1+\delta}{2d} \log \frac{\frac{1+\delta}{2d}}{1 - \delta/2d} + \frac{1-\delta}{2d} \log \frac{\frac{1-\delta}{2d}}{\frac{1+\delta}{2d}} = \frac{\delta}{d} \log \frac{1+\delta}{1-\delta} \leq \frac{3}{d}\delta^2, \quad \delta \leq \frac{1}{2}$$

So we conclude $M_n(\Theta, P) \geq \frac{d}{2} \cdot \frac{MD\delta}{d} \cdot \left(1 - \sqrt{\frac{1}{2} \cdot \frac{3k}{d}\delta^2}\right) = \frac{MD\delta}{2} \cdot \left(1 - \sqrt{\frac{3}{2d}k\delta^2}\right)$.

Setting $\delta^2 = \frac{d}{6k}$, which is possible if $\frac{d}{6} \leq \frac{1}{2} \cdot k$, $\sqrt{\frac{3}{2d}k\delta^2} = \frac{1}{2}$.

$$M_k(\Theta, P) \geq \frac{MD}{2} \cdot \sqrt{\frac{d}{6k}} \cdot \frac{1}{2} = \frac{MD}{4} \sqrt{\frac{d}{6k}} \qquad \text{whenever} \qquad k \geq \frac{d}{3}.$$

Let's say $[-D_{in}, D_{in}]^d \subseteq \Theta \subseteq [-D_{out}, D_{out}]^d$. Then,

$$\frac{D_{in} \cdot M}{4} \sqrt{\frac{d}{6k}} \leq M_k(\Theta, P) \leq \frac{3}{2} D_{out} \cdot M \cdot \sqrt{\frac{d}{k}} \qquad \text{from previous result.}$$