

Question 1: You toss a coin 10,000 times. How many heads do you see?

Question 2: Coupon collector problem. Have N coupons and we need to collect them all. How many coupons do we need to sample to get all N ?

Question 3: Largest common subsequence problem: have sequences X_1, \dots, X_n and Y_1, \dots, Y_n of iid Bern(1/2) random variables. What is the largest k such that there exist $i_1 < i_2 < \dots < i_k$ and $j_1 < j_2 < \dots < j_k$ such that $X_{i_1} = Y_{j_1}, \dots, X_{i_k} = Y_{j_k}$?

Question 1: we have various possible answers:

- 5,000. Indeed if we let X_i be the indicator of the event that we see heads on the i th toss, the number of heads is $S = \sum_{i=1}^{10000} X_i$ and $\mathbb{E}S = 5000$. But $\mathbb{P}(S = 5000) = \binom{10000}{5000} 2^{-10000} \approx 0.008$.
- Weak Law of Large Numbers: let $(X_i)_{i \geq 1}$ be iid with finite expectation μ and finite second moments. Then for every $\varepsilon > 0$,

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| > \varepsilon\right) \xrightarrow{n \rightarrow \infty} 0.$$

Therefore for large enough n , the number of heads lies in $[n(1/2 - \varepsilon), n(1/2 + \varepsilon)]$ with high probability. The main problem is that this is an asymptotic result - we don't know how large n should be.

- Central Limit Theorem: let $(X_i)_{i \geq 1}$ be iid with finite mean μ and finite second moment $\sigma^2 + \mu^2$. Then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) \xrightarrow{d} \mathcal{N}(0, 1).$$

Therefore $\sum_{i=1}^n (X_i - \mu)$ has deviations of the order $\sqrt{n}\sigma$. Suppose we pretend 10000 is big: then

$$\begin{aligned} S = \sum_{i=1}^{10000} X_i &\in [5000 - Q^{-1}(0.005)\sqrt{100}/2, 5000 + Q^{-1}(0.005)\sqrt{100}/2] \\ &\approx [5000 \pm 128] \end{aligned}$$

with probability 0.99, where $Q(x) = \mathbb{P}(Z \geq x)$ for $Z \sim \mathcal{N}(0, 1)$. However we have the same issue again - is $n = 10000$ large enough?

We can however give some non-asymptotic answers to Question 1:

Proposition (Chebyshev's inequality). Let X be any random variable with mean μ and variance σ^2 . Then

$$\mathbb{P}(|X - \mu| > t) \leq \frac{\sigma^2}{t^2}.$$

With this, we have

$$\mathbb{P}\left(\left|\sum_{i=1}^{10000} X_i - 5000\right| > t\right) \leq \frac{10000 \times \frac{1}{4}}{t^2} = \frac{2500}{t^2}.$$

So in particular, if $t = 500$ the RHS is 0.01. So we have $S \in [4500, 5500]$ with probability 0.99. However note that this is a weaker result than what the Central Limit Theorem gives.

Question 2: the number of samples S is equal to $\sum_{i=1}^N X_i$ where $X_i \sim \text{Geo}(i/N)$. Thus $\mathbb{E}S = \sum_{i=1}^N \frac{N}{i} = N \sum_{i=1}^N \frac{1}{i} \approx N \log N$.

Question 3: we have a function $f(X_1, \dots, X_n, Y_1, \dots, Y_n)$ which gives the longest common subsequence. It turns out this function is “smooth” in a certain sense, for which we can use “Talagrand’s Principle”.

Chernoff-Cr  mer method

Theorem (Markov's inequality). *Let Y be a non-negative random variable with finite expectation. Then for any $t > 0$ we have*

$$\mathbb{P}(Y \geq t) \leq \frac{\mathbb{E}Y}{t}.$$

Proof. Note $t\mathbb{1}(Y \geq t) \leq Y$ and integrate. \square

Corollary. *Let Y be a random variable. Suppose $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ is increasing and such that $\mathbb{E}|\phi(Y)| < \infty$. Then*

$$\mathbb{P}(Y \geq t) \leq \mathbb{P}(\phi(Y) \geq \phi(t)) \leq \frac{\mathbb{E}\phi(Y)}{\phi(t)}.$$

Note that for a random variable Z , letting $Y = |Z - \mathbb{E}Z|$ and $\phi : t \mapsto t^2$ gives Chebyshev's inequality $\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq \frac{\text{Var}(Z)}{t^2}$.

Could also take $\phi : t \mapsto t^q$ for any $q > 0$ to conclude $\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq \frac{\mathbb{E}|Z - \mathbb{E}|^q}{t^q}$.

Consider instead $\phi : t \mapsto e^{\lambda t}$ for $\lambda > 0$. Then we get

$$\mathbb{P}(Z \geq t) \leq \frac{\mathbb{E}[e^{\lambda Z}]}{e^{\lambda t}}.$$

Define $F(\lambda) = \mathbb{E}e^{\lambda Z}$, the *moment generating function* of Z . Define $\psi_Z(\lambda) = \log \mathbb{E}[e^{\lambda Z}]$. If X_1, \dots, X_n are independent and $Z = \sum_{i=1}^n X_i$ then it is clear that $\psi_Z(\lambda) = \sum_{i=1}^n \psi_{X_i}(\lambda)$. So we have

$$\mathbb{P}(Z \geq t) \leq \inf_{\lambda \geq 0} e^{\psi_Z(\lambda) - \lambda t}.$$

Now define $\psi_Z^*(t) = \sup_{\lambda \geq 0} (\lambda t - \psi_Z(\lambda))$ and write $\mathbb{P}(Z \geq t) \leq e^{-\psi_Z^*(t)}$. This is known as the *Chernoff bound*, and ψ_Z^* is known as the *Chernoff-Cr  mer transform*.

Properties of ψ_Z and ψ_Z^*

1. ψ_Z is convex and infinitely differentiable on $(0, b)$ where $b = \sup\{\lambda : \psi_Z(\lambda) < \infty\}$. Indeed

$$\begin{aligned} F(\theta x + (1 - \theta)y) &= \mathbb{E}[e^{\theta x Z} e^{(1 - \theta)y Z}] \\ &\leq \mathbb{E}[e^{x Z}]^\theta \mathbb{E}[e^{y Z}]^{1 - \theta}. \end{aligned} \quad (\text{H  lder with } 1/p = \theta, 1/q = 1 - \theta)$$

2. $\psi_Z^* \geq 0$ and it is convex (follows from the definition).

3. Suppose $t \geq \mathbb{E}Z$. Then $\psi_Z^*(t) = \sup_{\lambda} (\lambda t - \psi_Z(\lambda))$. Indeed we'll show $\lambda t - \psi_Z(\lambda) \leq 0$ whenever $\lambda < 0$. We have

$$\begin{aligned}\mathbb{E}[e^{\lambda Z}] &\geq e^{\lambda \mathbb{E}Z} && \text{(Jensen)} \\ \implies \psi_Z(\lambda) &\geq \lambda \mathbb{E}Z \\ \implies \lambda t - \psi_Z(\lambda) &\leq \lambda t - \lambda \mathbb{E}Z = \lambda(t - \mathbb{E}Z) \leq 0.\end{aligned}$$

Example. Let $Z \sim \mathcal{N}(0, v)$. We want to upper bound $\mathbb{P}(Z \geq t)$ for $t > 0$. We have

$$\begin{aligned}\mathbb{E}[e^{\lambda Z}] &= \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi v}} e^{-\frac{t^2}{2v}} e^{\lambda t} dt \\ &= \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi v}} e^{-\frac{(t-\lambda v)^2}{2v}} e^{\frac{v\lambda^2}{2}} dt \\ &= e^{\frac{v\lambda^2}{2}}.\end{aligned}$$

Hence $\psi_Z^*(t) = \sup_{\lambda} \left(\lambda t - \frac{\lambda^2 v}{2} \right)$ (for $t > 0 = \mathbb{E}Z$). Differentiating we see the optimal value is $\lambda = t/v$. Plugging this in gives $\psi_Z^*(t) = \frac{t^2}{2v}$. Thus

$$\mathbb{P}(Z \leq t) \leq e^{-\frac{t^2}{2v}}.$$

Sub-Gaussian random variables

Definition. A random variable Y with $\mathbb{E}Y = 0$ is *sub-Gaussian* with variance parameter v if

$$\psi_Y(\lambda) < \frac{\lambda^2 v}{2} \quad \forall \lambda \in \mathbb{R}.$$

The set of sub-Gaussian random variables with variance parameter v is denoted $\mathcal{G}(v)$.

1. It is clear from the above that if $Y \in \mathcal{G}(v)$ then $\mathbb{P}(Y \geq t) \leq e^{-t^2/2v}$ and $\mathbb{P}(Y \leq -t) \leq e^{-t^2/2v}$.
2. If $Y_i \in \mathcal{G}(v_i)$ for $i = 1, \dots, n$ are independent then $\sum_{i=1}^n Y_i \in \mathcal{G}(\sum_{i=1}^n v_i)$ (immediate by additivity of $\psi(\cdot)$).
3. If $Y \in \mathcal{G}(v)$ then $\text{Var}(Y) \leq v$ (see Example Sheet).

Theorem. The following are equivalent for suitable v, b, c, d

1. $Y \in \mathcal{G}(v)$;
2. $\max\{\mathbb{P}(Y \geq t), \mathbb{P}(Y \leq -t)\} \leq e^{-\frac{t^2}{2b}}$ for all $t > 0$;
3. $\mathbb{E}Y^{2q} \leq q!c^q$ for all $q \geq 1$;
4. $\mathbb{E}[e^{dY^2}] \leq 2$.

Proof. Not given. □

Lemma (Hoeffding's lemma). Let Y be supported on $[a, b]$ and suppose $\mathbb{E}Y = 0$. Then $\psi_Y''(\lambda) \leq \frac{(b-a)^2}{4}$, and so $Y \in \mathcal{G}\left(\frac{(b-a)^2}{4}\right)$.

Proof. We have

$$\psi'_Y(\lambda) = \frac{\mathbb{E}[Y e^{\lambda Y}]}{\mathbb{E}[e^{\lambda Y}]} \implies \psi''_Y(\lambda) = \frac{\mathbb{E}[e^{\lambda Y}] \mathbb{E}[Y^2 e^{\lambda Y}] - (\mathbb{E}[Y e^{\lambda Y}])^2}{\mathbb{E}[e^{\lambda Y}]^2}.$$

So

$$\begin{aligned} \psi''_Y(\lambda) &= \int_{\mathbb{R}} y^2 \underbrace{\frac{e^{\lambda y}}{\mathbb{E}[e^{\lambda Y}]}}_{:=dQ(y)} d\mu_Y(y) - \left(\int_{\mathbb{R}} y \frac{e^{\lambda y}}{\mathbb{E}[e^{\lambda Y}]} d\mu_Y(y) \right)^2 \\ &= \text{Var}_{Y \sim Q}(Y) \geq 0 \end{aligned}$$

noting that Q is supported on $[a, b]$. If $Y \in [a, b]$ almost-surely then note

$$\text{Var}(Y) = \text{Var}\left(Y - \frac{a+b}{2}\right) \leq \mathbb{E}\left[\left(Y - \frac{a+b}{2}\right)^2\right] \leq \frac{(b-a)^2}{4}.$$

To finish, observe that $\psi_Y(\lambda) = \psi_Y(0) + \lambda \psi'_Y(0) + \frac{\lambda^2}{2} \psi''_Y(\theta)$ for some $\theta \in [0, \lambda]$. Thus $\psi_Y(\lambda) \leq \frac{\lambda^2(b-a)^2}{8}$. \square

Theorem (Hoeffding's inequality). *Let Y_1, \dots, Y_n be independent random variables with Y_i having support on $[a_i, b_i]$. Then*

$$\mathbb{P}\left(\sum_{i=1}^n (Y_i - \mathbb{E}Y_i) \geq t\right) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

Proof. Trivial by Hoeffding's lemma and additivity of the variance parameters. \square

Theorem (Bennett's inequality). *For $1 \leq i \leq n$, let X_i be independent random variables satisfying $\mathbb{E}X_i = 0$, $\text{Var}(X_i) = \sigma_i^2$ and let $v = \sum_{i=1}^n \sigma_i^2$. Further assume the X_i are bounded by some $C > 0$ almost-surely. Then*

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{v}{C^2} h_1\left(\frac{Ct}{v}\right)\right)$$

where $h_1(x) = (1+x) \log(1+x) - x$ for $x > 0$. Furthermore, using the inequality $h_1(x) \geq \frac{x^2}{2(1+x/3)}$ we obtain

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{t^2}{2(v + Ct/3)}\right).$$

Example. Suppose $X_i \sim \text{Bern}(p_n)$ are independent for $1 \leq i \leq n$. Then

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{2t^2}{n}\right) \quad (\text{Hoeffding})$$

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{t^2}{np_n(1-p_n) + t/3}\right). \quad (\text{Bennett})$$

Note that if $p_n \ll q$, e.g. $p_n = 1/\sqrt{n}$, Hoeffding will stay the same, i.e. of order $e^{-\frac{2t^2}{n}}$ (only depends on support, not variance). However, Bennet will be of the order $e^{-\frac{t^2}{\sqrt{n+t/3}}}$.

Proof. We have

$$\begin{aligned}
 \mathbb{E}[e^{\lambda X_i}] &= \sum_{k \geq 0} \frac{\lambda^k}{k!} \mathbb{E}[X_i^k] \\
 &\leq 1 + \sum_{k \geq 2} \frac{\lambda^k}{k!} \mathbb{E}[C^{k-2} X_i^2] \\
 &= 1 + \sum_{k \geq 2} \frac{\lambda^k C^{k-2} \sigma_i^2}{k!} \\
 &= 1 + \frac{\sigma_i^2}{C^2} (e^{\lambda C} - \lambda C - 1) \\
 &\leq \exp \left(\frac{\sigma_i^2}{C^2} (e^{\lambda C} - \lambda C - 1) \right). \quad ((1+x) \leq e^x)
 \end{aligned}$$

This implies

$$\mathbb{E}^{\lambda S} \leq \exp \left(\frac{v}{C^2} (e^{\lambda C} - \lambda C - 1) \right)$$

and so

$$\psi_S(\lambda) \leq \underbrace{\frac{v}{C^2} (e^{\lambda C} - \lambda C - 1)}_{:= \tilde{\psi}(\lambda)}.$$

This means that

$$\psi_S^*(t) \geq \tilde{\psi}^*(t)$$

and

$$\mathbb{P}(S \geq t) \leq \exp(-\psi_S^*(t)) \leq \exp(-\tilde{\psi}^*(t)) = \exp \left(-\frac{v}{C^2} h_1 \left(\frac{Ct}{v} \right) \right)$$

where the last equality is by a result from Example Sheet 1. \square

Efron-Stein Inequality

We want to bound $\text{Var}(Z)$ where $Z = f(X_1, \dots, X_n)$ for independent X_i 's (or even just uncorrelated). If $Z - \mathbb{E}Z = \sum_{i=1}^n \Delta_i$ for $\Delta_1, \dots, \Delta_n$ uncorrelated and with 0 mean we have $\text{Var}(Z) = \sum_{i=1}^n \mathbb{E}[\Delta_i^2]$. Define $\mathbb{E}_i Z = \mathbb{E}[Z|X_{1:i}]^1$ where $X_{1:i} = (X_1, \dots, X_i)$.

Set $\Delta_i = \mathbb{E}_i Z - \mathbb{E}_{i-1} Z$. Then $Z - \mathbb{E}Z = \sum_{i=1}^n \Delta_i$. Also $\mathbb{E}\Delta_i = 0$ by the tower property of conditional expectation. Suppose $i < j$ so

$$\begin{aligned}\mathbb{E}[\Delta_i \Delta_j] &= \mathbb{E}[\mathbb{E}[\Delta_i \Delta_j | X_{1:i}]] \\ &= \mathbb{E}[\Delta_i \mathbb{E}[\Delta_j | X_{1:i}]].\end{aligned}$$

Note that $\mathbb{E}[\Delta_j | X_{1:i}] = \mathbb{E}[\mathbb{E}_j Z | X_{1:i}] - \mathbb{E}[\mathbb{E}_{j-1} Z | X_{1:i}] = \mathbb{E}_i Z - \mathbb{E}_{i-1} Z = 0$. Thus $\mathbb{E}[\Delta_i \Delta_j] = 0$ and so the Δ_i 's are uncorrelated.

Thus $\text{Var}(Z) = \sum_{i=1}^n \mathbb{E}[\Delta_i^2]$ regardless of the correlation between the X_i (though we still assume independence of the X_i going forward).

Define $\mathbb{E}^{(i)} Z = \mathbb{E}[Z | X_{1:i-1}, X_{i+1:n}]$. Then $\Delta_i = \mathbb{E}_i Z - \mathbb{E}_{i-1} Z = \mathbb{E}_i (Z - \mathbb{E}^{(i)} Z)$. Indeed we have $\mathbb{E}_i[\mathbb{E}^{(i)} Z] = \mathbb{E}[\mathbb{E}[Z | X^{(i)}] | X_{1:i}] = \mathbb{E}[\mathbb{E}[Z | X^{(i)}] | X_{1:i-1}]$ by independence and $\mathbb{E}[\mathbb{E}[Z | X^{(i)}] | X_{1:i-1}] = \mathbb{E}[Z | X_{1:i-1}]$ since $\sigma(X_{1:i-1}) \subseteq \sigma(X^{(i)})$.

Therefore

$$\Delta_i^2 = (\mathbb{E}_i (Z - \mathbb{E}^{(i)} Z))^2 \leq \mathbb{E}_i [(Z - \mathbb{E}^{(i)} Z)^2]$$

almost-surely by conditional Jensen.

Hence we have

$$\begin{aligned}\text{Var}(Z) &= \sum_{i=1}^n \mathbb{E}[\Delta_i^2] \\ &\leq \sum_{i=1}^n \mathbb{E}[(Z - \mathbb{E}^{(i)} Z)^2] \\ &= \sum_{i=1}^n \mathbb{E}[\mathbb{E}[(Z - \mathbb{E}^{(i)} Z)^2] | X^{(i)}] \\ &= \mathbb{E} \left[\sum_{i=1}^n \text{Var}^{(i)}(Z) \right].\end{aligned}$$

This is called the *Efron-Stein inequality*.

¹For a rigorous definition of this conditional expectation see Part III Advanced Probability

To summarise:

Theorem (Efron-Stein Inequality). *Let X_1, \dots, X_n be independent random variables and let $Z = f(X_1, \dots, X_n)$ be a square integrable function of $X = X_{1:n}$. Then*

$$\text{Var}(Z) \leq \sum_{i=1}^n \mathbb{E}[(Z - \mathbb{E}^{(i)} Z)^2] = \underbrace{\sum_{i=1}^n \text{Var}^{(i)}(Z)}_{:=v}.$$

Proposition. Define X'_1, \dots, X'_n to be independent copies of X_1, \dots, X_n respectively. Set $Z'_i = f(X^{(i)}, X'_i)$. Then

$$v = \sum_{i=1}^n \mathbb{E}[(Z - Z'_i)_+^2] = \sum_{i=1}^n \mathbb{E}[(Z - Z'_i)_-^2] = \frac{1}{2} \sum_{i=1}^n \mathbb{E}[(Z - Z'_i)^2].$$

Also

$$v = \inf_{Z_1, \dots, Z_n} \sum_{i=1}^n \mathbb{E}[(Z - Z_i)^2]$$

where Z_i is some function of $X^{(i)}$.

Proof. Note that if X, Y are iid then

$$\text{Var}(X) = \frac{1}{2} \mathbb{E}[(X - Y)^2] = \mathbb{E}[(X - Y)_+^2] = \mathbb{E}[(X - Y)_-^2]$$

since $(X - Y)_+, (X - Y)_-$ have the same distribution. For the final expression, note $\text{Var}(X) = \inf_a \mathbb{E}[(X - a)^2]$. Then $\text{Var}^{(i)}(Z) = \inf_{Z_i} \mathbb{E}[(Z - Z_i)^2 | X^{(i)}]$ where Z_i is $X^{(i)}$ -measurable. \square

Functions with bounded-differences property

We say f satisfies the *bounded differences property* with constants c_1, \dots, c_n if

$$\sup_{x_1, \dots, x_n, x'_i} |f(x_1, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i.$$

If $Z = f(X_1, \dots, X_n)$ where the X_i are independent and f satisfying bounded differences, we'll show that $\text{Var}(Z) \leq \sum_{i=1}^n \frac{c_i^2}{4}$. To see this, set

$$Z_i = \frac{1}{2} \left(\inf_{x_i} f(X^{(i)}, x_i) + \sup_{x_i} f(X^{(i)}, x_i) \right).$$

Then

$$v \leq \sum_{i=1}^n \mathbb{E}[(Z - Z_i)^2] \leq \sum_{i=1}^n \frac{c_i^2}{4}.$$

Example. Let X_1, \dots, X_n be independent and supported on $[0, 1]$. Define $f(X_{1:n})$ to be the smallest number of size 1 bins needed to “pack” X_1, \dots, X_n . Note f satisfies the bounded differences property with $c_i = 1$ for all i . Therefore $\text{Var}(Z) \leq \frac{n}{4}$. Suppose now the X_i are iid uniform on $[0, 1]$. Then $\mathbb{E}f(X_1, \dots, X_n) \approx Cn$ while the standard deviation is of order at most \sqrt{n} , giving tight confidence intervals for large n .

Example. Let $X_1, \dots, X_n, Y_1, \dots, Y_n$ be iid Bernoulli with parameter $1/2$. Let $f(X_{1:n}, Y_{1:n})$ be the longest common subsequence between $X_{1:n}$ and $Y_{1:n}$. Then f satisfies bounded differences with $c_i = 1$ for all i . Thus $\text{Var}(Z) \leq n/2$. It is known that $\mathbb{E}[Z] \sim [0.75n, 0.837n]$. So again Z is very concentrated about its mean for large n .

Example. The chromatic number $\chi(G)$ of a graph G is the smallest number of colours needed to colour vertices of G such that no two neighbouring vertices have the same colour. Let X_{ij} be iid Bernoulli of parameter p for $1 \leq i < j \leq n$. We construct a random graph G on vertex set $\{1, \dots, n\}$ by saying $\{i, j\} \in E$ iff $X_{ij} = 1$. Take f such that $f(\{X_{ij}\}_{1 \leq i < j \leq n}) = \chi(G)$. Then f again satisfies bounded differences with $c_{ij} = 1$ for all $1 \leq i < j \leq n$. Hence $\text{Var}(\chi(G)) \leq \frac{1}{4} \binom{n}{2}$. It is known that $\mathbb{E}[\chi(G)] \approx n/\log n$. This gives a poor confidence interval.

However, we can fix this bound by considering $Y_i = (X_{1,i+1}, \dots, X_{i,i+1})$. Observe that Y_1, \dots, Y_{n-1} are independent and $\chi(G)$ is some function \hat{f} of Y_1, \dots, Y_{n-1} . It can be shown that we still have bounded differences with $c_1 = \dots = c_{n-1} = 1$. This gives $\text{Var}(\chi(G)) \leq \frac{n-1}{4}$ and thus we have a good confidence interval now.

Theorem (Convex Poincaré Inequality). *Let X_1, \dots, X_n be independent and supported on $[0, 1]$. Let f be a separately convex function (i.e convex in each variable) over $[0, 1]^n$ which has partial derivatives. Then*

$$\text{Var}(f(X)) \leq \mathbb{E}[\|\nabla f(X)\|^2].$$

Remark. Jointly convex functions are separately convex so this inequality holds for such functions too.

Proof. We have

$$\text{Var}(Z) \leq \sum_{i=1}^n \mathbb{E}[(Z - Z_i)^2]$$

where Z_i is $X^{(i)}$ -measurable. Let $Z_i = \inf_x f(X^{(i)}, x)$. Then

$$Z - Z_i = f(X_1, \dots, X_n) - f(X_1, \dots, X_{i-1}, x^*, x_{i+1}, \dots, X_n) = f(X^{(i)}, X_i) - f(X^{(i)}, x^*) \geq 0$$

where x^* achieves the infimum of $f(X^{(i)}, x)$ over x . If g is convex then $g(y) \geq g(x) + g'(x)(y - x)$. Hence

$$f(X^{(i)}, X_i) - f(X^{(i)}, x^*) \leq \frac{\partial f}{\partial x_i}(X) \cdot (x^* - X_i).$$

Squaring gives

$$(Z - Z_i)^2 \leq \left[\frac{\partial f}{\partial x_i}(X)(x^* - X_i) \right]^2 \leq \left[\frac{\partial f}{\partial x_i}(X) \right]^2.$$

□

Example. Let $X \in \mathbb{R}^{n \times d}$ with $\mathbb{E}X_{ij} = 0$ for all i, j and with all entries independent and supported on $[-1, 1]$. Let

$$\sigma_1(X) = \max_{\|v\|_2=1} \|Xv\|_1 = \max_{\|U\|_2=1, \|v\|_2=1} U^T X v.$$

Can show the triangle inequality holds so

$$|\sigma_1(A) - \sigma_1(B)| \leq \sigma_1(A - B).$$

Can also show by Cauchy-Schwarz that

$$\sigma_1(A)^2 \leq \sum_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}} A_{ij}^2 = \|A\|_F^2.$$

Thus

$$|\sigma_1(A) - \sigma_1(B)| \leq \|A - B\|_F.$$

Therefore σ_1 is Frobenius-1-Lipschitz. This means (assuming derivatives exist) $\|\nabla \sigma_1(X)\| \leq 1$. So using the convex Poincaré inequality, $\text{Var}(\sigma_1(X)) \leq 4$.

Theorem (Gaussian Poincaré inequality). *Let X_1, \dots, X_n be iid $\mathcal{N}(0, 1)$. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuously differentiable. Then $\text{Var}(f(X)) \leq \mathbb{E}[\|\nabla f(X)\|^2]$.*

Proof. It is enough to show the $n = 1$ case. Indeed if the $n = 1$ case is true we have

$$\text{Var}(f(X_1, \dots, X_n)) \leq \sum_{i=1}^n \mathbb{E}[\text{Var}^{(i)}(Z)]$$

by Efron-Stein. Also

$$\text{Var}^{(i)}(Z) = \mathbb{E}[(Z - \mathbb{E}^{(i)} Z)^2 | X^{(i)}] \leq \mathbb{E} \left[\left(\frac{\partial f}{\partial x_i}(X) \right)^2 | X^{(i)} \right]$$

by the $n = 1$ case, so we get the general case.

Now we prove the $n = 1$ case. Let X_1, \dots, X_n be iid (Rademacher) symmetric $\text{Ber}(1/2)$ (i.e takes values ± 1 with probabilities $1/2$). Define $S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i$ so $S_n \xrightarrow{d} \mathcal{N}(0, 1)$ by the CLT. Then

$$\begin{aligned} \text{Var}(f(S_n)) &\leq \sum_{i=1}^n \mathbb{E}[\text{Var}^{(i)}(f(S_n))] \\ &= \sum_{i=1}^n \mathbb{E} \left[\frac{1}{4} \left(f(S_n - \frac{X_i}{\sqrt{n}} + \frac{1}{\sqrt{n}}) - f(S_n - \frac{X_i}{\sqrt{n}} - \frac{1}{\sqrt{n}}) \right)^2 \right]. \end{aligned}$$

For the rest of the proof we assume f is twice-continuously differentiable on a bounded domain. Then

$$\begin{aligned} f(S_n - \frac{X_i}{\sqrt{n}} + \frac{1}{\sqrt{n}}) &= f(S_n) + f'(S_n) \frac{1 - X_i}{\sqrt{n}} + f''(\theta_1) \frac{(1 - X_i)^2}{2n} \\ f(S_n - \frac{X_i}{\sqrt{n}} - \frac{1}{\sqrt{n}}) &= f(S_n) - f'(S_n) \frac{1 + X_i}{\sqrt{n}} + f''(\theta_2) \frac{(1 + X_i)^2}{2n} \end{aligned}$$

so

$$|f(S_n - \frac{X_i}{\sqrt{n}} + \frac{1}{\sqrt{n}}) - f(S_n - \frac{X_i}{\sqrt{n}} - \frac{1}{\sqrt{n}})| \leq |f'(S_n)| \frac{2}{\sqrt{n}} + \|f''\|_{\infty} \frac{2}{n}.$$

Hence

$$\begin{aligned} &|f(S_n - \frac{X_i}{\sqrt{n}} + \frac{1}{\sqrt{n}}) - f(S_n - \frac{X_i}{\sqrt{n}} - \frac{1}{\sqrt{n}})|^2 \\ &\leq |f'(S_n)|^2 \frac{4}{n} + \frac{4\|f''\|_{\infty}^2}{n^2} + \frac{8|f'(S_n)|\|f''\|_{\infty}}{n^{3/2}} \end{aligned}$$

and summing over $\{1, \dots, n\}$ we get

$$\text{Var}(f(S_n)) \leq \mathbb{E}[f'(S_n)^2] + \frac{\|f''\|_{\infty}^2}{n} + \frac{8\mathbb{E}[|f'(S_n)|]\|f''\|_{\infty}}{n^{1/2}}$$

so taking $n \rightarrow \infty$ gives the result. \square

Entropy

Definition. For a random variable taking values on a discrete set \mathcal{X} with PMF P_X , the *Shannon entropy* is defined as $H(X) = H(P_X) = \mathbb{E}[-\log P_X(X)]$.

Definition. Given two probability measures P, Q on a discrete set \mathcal{X} , define the *relative entropy* or *Kullback-Leibler divergence* $D(Q\|P) = \sum_x q(x) \log \frac{q(x)}{p(x)}$ where p, q are the PMF's of P, Q respectively.

Some basic properties of relative entropy are

1. $D(Q\|P) \geq 0$ with equality iff $Q = P$;
2. $D(Q\|P)$ is jointly convex, i.e

$$D(\lambda Q_1 + (1 - \lambda)Q_2\|\lambda P_1 + (1 - \lambda)P_2) \leq \lambda D(Q_1\|P_1) + (1 - \lambda)D(Q_2\|P_2).$$

Suppose $|\mathcal{X}| < \infty$, then

$$D(Q\|U) = \log |\mathcal{X}| - H(Q)$$

where $U \sim \text{Uniform}(\mathcal{X})$.

Definition. We define the *conditional entropy* $H(Y|X)$ by

$$\begin{aligned} H(Y|X) &= \mathbb{E}[-\log P_{Y|X}(Y|X)] \\ &= - \sum_{x,y} P_{X,Y}(x,y) \log P_{Y|X}(y|x) \\ &= \sum_x H(Y|X=x) P_X(x) = \sum_x H(P_{Y|X=x}) P_X(x). \end{aligned}$$

Note $H(Y|X) \leq H(Y)$ by concavity of H together with Jensen. We define the *joint entropy* $H(X, Y)$ by

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y) = \mathbb{E}[-\log P_{X,Y}(X, Y)].$$

Theorem (Chain rule). $H(X_1, \dots, X_n) = \sum_{i=1}^n H(X_i|X_{1:i-1})$.

Proof. We have

$$\begin{aligned} H(X_1, \dots, X_n) &= \mathbb{E}[-\log P_{X_{1:n}}(X_{1:n})] \\ &= \mathbb{E} \left[-\log \prod_{i=1}^n P_{X_i|X_{1:i-1}}(X_i|X_{1:i-1}) \right] \\ &= \sum_{i=1}^n \mathbb{E}[-\log P_{X_i|X_{1:i-1}}(X_i|X_{1:i-1})] \\ &= \sum_{i=1}^n H(X_i|X_{1:i-1}). \end{aligned}$$

□

Theorem (Chain rule for KL-divergence). *Let P, Q be measures on \mathcal{X}^n . Then*

$$D(Q\|P) = D(Q\|P) = \sum_{i=1}^n D(Q_{X_i|X_{1:i-1}}\|P_{X_i|X_{1:i-1}}|Q_{X_{1:i-1}}).$$

Proof. We have

$$\begin{aligned} D(Q\|P) &= \sum_{x_{1:n}} q(x_{1:n}) \log \frac{q(x_{1:n})}{p(x_{1:n})} \\ &= \mathbb{E}_Q \left[\log \frac{q(X_{1:n})}{p(X_{1:n})} \right] \\ &= \mathbb{E}_Q \left[\log \prod_{i=1}^n \frac{q(X_i|X_{1:i-1})}{p(X_i|X_{1:i-1})} \right] \\ &= \sum_{i=1}^n \mathbb{E}_Q \left[\log \frac{q(X_i|X_{1:i-1})}{p(X_i|X_{1:i-1})} \right]. \end{aligned}$$

Note that

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}_Q \left[\log \frac{q(X_i|X_{1:i-1})}{p(X_i|X_{1:i-1})} \right] &= \sum_{x_{1:i}} q(x_{1:i}) \log \frac{q(x_i|x_{1:i-1})}{p(x_i|x_{1:i-1})} \\ &= \sum_{x_{1:i-1}} q(x_{1:i-1}) \left[\sum_{x_i} q(x_i|x_{1:i-1}) \log \frac{q(x_i|x_{1:i-1})}{p(x_i|x_{1:i-1})} \right] \\ &= \mathbb{E}_{Q_{X_{1:i-1}}} [D(Q_{X_i|X_{1:i-1}}\|P_{X_i|X_{1:i-1}})] \\ &:= D(Q_{X_i|X_{1:i-1}}\|P_{X_i|X_{1:i-1}}|Q_{X_{1:i-1}}). \end{aligned}$$

Hence

$$D(Q\|P) = \sum_{i=1}^n D(Q_{X_i|X_{1:i-1}}\|P_{X_i|X_{1:i-1}}|Q_{X_{1:i-1}}).$$

□

Usually we'll have $P = P_1 \otimes P_2 \otimes \dots \otimes P_n$, which simplifies this expression. If $Q = Q_1 \otimes Q_2 \otimes \dots \otimes Q_n$ then it simplifies further to

$$D(Q\|P) = \sum_{i=1}^n D(Q_i\|P_i).$$

Theorem (Han's inequality for Shannon entropy). *We have*

$$H(X_{1:n}) \leq \frac{\sum_{i=1}^n H(X^{(i)})}{n-1}.$$

Example. Let $X_{1:n}$ be sampled iid from the uniform distribution on $A \subseteq \mathbb{Z}^n$. Then $H(X_{1:n}) = \log |A|$. Then Han's inequality implies

$$\log |A| \leq \frac{\log |A^{(i)}|}{n-1} \implies |A| \leq \left(\prod_{i=1}^n |A^{(i)}| \right)^{1/(n-1)}$$

which is called the Loomis-Whitney inequality.

Lemma. *We have*

$$H(X|Y, Z) \leq H(X|Y).$$

Proof. We have

$$\begin{aligned} H(X|Y, Z) &= \sum_{y,z} H(P_{X|Y=y, Z=z}) P_{YZ}(y, z) \\ &= \sum_y P_Y(y) \left[\sum_z P_{Z|Y}(z|y) H(P_{X|Y=y, Z=z}) \right] \\ &\leq \sum_y P_Y(y) H \left(\sum_z P_{Z|Y}(z|y) P_{X|Y=y, Z=z} \right) \quad (\text{concavity of } H) \\ &= \sum_y P_Y(y) H(P_{X|Y=y}) \\ &= H(X|Y). \end{aligned}$$

□

Now we prove:

Theorem (Han's inequality for Shannon entropy). *We have*

$$H(X_{1:n}) \leq \frac{\sum_{i=1}^n H(X^{(i)})}{n-1}.$$

Proof. We have

$$\begin{aligned} H(X_{1:n}) &= H(X^{(i)}) + H(X_i|X^{(i)}) \\ &\leq H(X^{(i)}) + H(X_i|X_{1:i-1}) \end{aligned}$$

by the previous lemma. Now summing over i and applying the chain rule gives

$$nH(X_{1:n}) \leq \sum_{i=1}^n H(X^{(i)}) + H(X_{1:n}).$$

□

Theorem (Han's inequality for KL-divergence). *Let \mathcal{X} be a countable set and let P, Q be measures on \mathcal{X}^n where $P = P_1 \otimes P_2 \otimes \dots \otimes P_n$. Then*

$$D(Q\|P) \geq \frac{1}{n-1} \sum_{i=1}^n D(Q_{X^{(i)}}\|P_{X^{(i)}})$$

or equivalently,

$$D(Q\|P) \leq \sum_{i=1}^n D(Q_{X_i|X^{(i)}}\|P_{X_i|Q_{X^{(i)}}}).$$

Remark.

$$D(Q\|P) = D(Q_{X^{(i)}}\|P_{X^{(i)}}) + D(Q_{X_i|X^{(i)}}\|\underbrace{P_{X_i|X^{(i)}}}_{P_{X_i}}|Q_{X^{(i)}})$$

Remark. If \mathcal{X} is finite and P_1, \dots, P_n are uniform over \mathcal{X} then we get Han's inequality for Shannon entropy.

Lemma. Let P, Q be measures on a discrete set $\mathcal{X} \times \mathcal{Y} \times \mathcal{Y}$. Then

$$D(Q_{Y|XZ}\|P_Y|Q_{XZ}) \geq D(Q_{Y|X}\|P_Y|Q_X).$$

Proof. We have

$$\begin{aligned} D(Q_{Y|XZ}\|P_Y|Q_{XZ}) &= \sum_{x,z} Q_{XZ}(x,z) D(Q_{Y|X=x,Z=z}\|P_Y) \\ &= \sum_x Q_X(x) \left[\sum_z Q_{Z|X}(z|x) D(Q_{Y|X=x,Z=z}\|P_Y) \right] \\ &= \sum_z Q_X(x) D \left(\sum_z Q_{Z|X}(z|x) Q_{Y|X=x,Z=z} \| P_Y \right) \\ &\quad \text{(convexity of } D) \\ &= \sum_z Q_X(x) D(Q_{Y|X=x}\|P_Y) \\ &= D(Q_{Y|X}\|P_Y|Q_X). \end{aligned}$$

□

Proof of Han's inequality for KL-divergence. We have

$$\begin{aligned} D(Q\|P) &= D(Q_{X^{(i)}}\|P_{X^{(i)}}) + D(Q_{X_i|X^{(i)}}\|P_{X_i}|Q_{X^{(i)}}) \\ &\geq D(Q_{X^{(i)}}\|P_{X^{(i)}}) + D(Q_{X_i|X_{1:i-1}}\|P_{X_i}|Q_{X_{1:i-1}}) \end{aligned}$$

by the previous lemma, so summing over i gives

$$nD(Q\|P) \geq \sum_{i=1}^n D(Q_{X^{(i)}}\|P_{X^{(i)}}) + D(Q\|P).$$

□

We have that $\text{Var}(Z) = \mathbb{E}Z^2 - [\mathbb{E}Z]^2 = \mathbb{E}[\phi(Z)] - \phi(\mathbb{E}Z)$ where $\phi(x) = x^2$.

Define $\text{Ent}(Z) := \mathbb{E}[Z \log Z] - \mathbb{E}[Z] \log \mathbb{E}[Z]$ for $Z \geq 0$, i.e take $\phi(x) = x \log x$. Since ϕ is convex, $\text{Ent}(Z) \geq 0$.

Suppose $Z = \frac{Q(X)}{P(X)}$ where $X \sim P$. Then $\mathbb{E}Z = 1$. Also

$$\text{Ent}(Z) = \mathbb{E} \left[\frac{Q(X)}{P(X)} \log \frac{Q(X)}{P(X)} \right] - 1 \log 1 = D(Q\|P).$$

Theorem (Han's inequality for Ent/Tensorisation of Ent). *Let X_1, \dots, X_n be independent random variables over \mathcal{X} (not necessarily discrete) and let $f : \mathcal{X}^n \rightarrow [0, \infty)$. Let $Z = f(X_{1:n})$. Then*

$$\text{Ent}(Z) \leq \sum_{i=1}^n \mathbb{E}[\text{Ent}^{(i)}(Z)]$$

where $\text{Ent}^{(i)}(Z) = \mathbb{E}^{(i)}[Z \log Z] - \mathbb{E}^{(i)} Z \log \mathbb{E}^{(i)} Z$, $\mathbb{E}^{(i)} Z = \mathbb{E}[Z | X^{(i)}]$.

Proof sketch. The case $Z \equiv 0$ is trivial so assume $Z \not\equiv 0$. WLOG we may also assume $\mathbb{E}Z = 1$. It is easy to check $\text{Ent}(aZ) = a\text{Ent}(Z)$ for $a > 0$. Since $\mathbb{E}Z = 1$ we have

$$\sum_{x_1, \dots, x_n} f(x_1, \dots, x_n) P_{X_{1:n}}(x_1, \dots, x_n) = 1.$$

Define $q(x_1, \dots, x_n) = f(x_{1:n}) P_{X_{1:n}}(x_{1:n})$. Then $\text{Ent}(Z) = D(Q \| P)$. By Han's inequality for KL-divergence we have

$$\text{Ent}(Z) = D(Q \| P) \leq \sum_{i=1}^n D(Q_{X_i | X^{(i)}} \| P_{X_i} | Q_{X^{(i)}})$$

and by a result on the examples sheet we have $D(Q_{X_i | X^{(i)}} \| P_{X_i} | Q_{X^{(i)}}) = \mathbb{E}[\text{Ent}^{(i)}(Z)]$. \square

Theorem (Herbst's argument). *Let Z be an integrable random variable such that for some $v > 0$ we have*

$$\text{Ent}(e^{\lambda Z}) \leq \frac{\lambda^2 v}{2} \mathbb{E}[e^{\lambda Z}]$$

for all $\lambda > 0$. Then $\Psi_{Z - \mathbb{E}Z}(\lambda) = \log \mathbb{E}[e^{\lambda(Z - \mathbb{E}Z)}] \leq \frac{\lambda^2 v}{2}$ for all $\lambda > 0$.

Proof. We have

$$\begin{aligned} \Psi_{Z - \mathbb{E}Z}(\lambda) &= \log \mathbb{E}[e^{\lambda Z}] - \lambda \mathbb{E}Z \\ \Psi'_{Z - \mathbb{E}Z}(\lambda) &= \frac{\mathbb{E}[Z e^{\lambda Z}]}{\mathbb{E}[e^{\lambda Z}]} - \mathbb{E}Z \\ \text{Ent}(e^{\lambda Z}) &= \mathbb{E}[e^{\lambda Z} \lambda Z] - \mathbb{E}[e^{\lambda Z}] \log \mathbb{E}[e^{\lambda Z}] = \mathbb{E}[e^{\lambda Z}] (\lambda \Psi'(\lambda) - \Psi(\lambda)). \end{aligned}$$

We have

$$\frac{\text{Ent}(e^{\lambda Z})}{\mathbb{E}[e^{\lambda Z}]} = \lambda \Psi'(\lambda) - \Psi(\lambda) \leq \frac{\lambda^2 v}{2} \text{ for } \lambda > 0.$$

Thus $\frac{\Psi'(\lambda)}{\lambda} - \frac{\Psi(\lambda)}{\lambda^2} = \left(\frac{\Psi(\lambda)}{\lambda} \right)' \leq v/2$. Hence integrating from 0 to λ gives

$$\frac{\Psi(\lambda)}{\lambda} \leq \frac{\lambda v}{2} \implies \Psi(\lambda) \leq \frac{\lambda^2 v}{2}.$$

\square

Theorem. *Let $f : \mathcal{X}^n \rightarrow \mathbb{R}$ satisfy the bounded differences property with c_1, \dots, c_n . Let X_1, \dots, X_n be independent and $Z = f(X_{1:n})$. Then for $t \geq 0$,*

$$\begin{aligned} \mathbb{P}(Z - \mathbb{E}Z > t) &\leq e^{-\frac{t^2}{2v}} \text{ where } v = \frac{\sum_{i=1}^n c_i^2}{4} \text{ and} \\ \mathbb{P}(Z - \mathbb{E}Z < -t) &\leq e^{-\frac{t^2}{2v}}. \end{aligned}$$

Proof. By tensorisation

$$\text{Ent}(e^{\lambda Z}) \leq \mathbb{E} \left[\sum_{i=1}^n \text{Ent}^{(i)}(e^{\lambda Z}) \right].$$

Assume the following lemma for now.

Lemma. *Let Y be bounded on $[a, b]$. Then*

$$\text{Ent}(e^{\lambda Y}) \leq \mathbb{E}(e^{\lambda Y}) \frac{(b-a)^2 \lambda^2}{8}.$$

Supposing the lemma is true we get

$$\text{Ent}^{(i)}(e^{\lambda Z}) \leq \mathbb{E}^{(i)}(e^{\lambda Z}) \frac{c_i^2 \lambda^2}{8}$$

and so

$$\text{Ent}(e^{\lambda Z}) \leq \mathbb{E}[e^{\lambda Z}] \frac{\lambda^2 v}{2}$$

then Herbst's argument shows $\Psi_{Z-\mathbb{E}Z}(\lambda) \leq \frac{\lambda^2 v}{2}$ and we use the Chernoff bound to get the result. \square

Proof of lemma. Recall that

$$\frac{\text{Ent}(e^{\lambda Y})}{\mathbb{E}[e^{\lambda Y}]} = \lambda \Psi'(\lambda) - \Psi(\lambda) = \int_0^\lambda t \psi''(t) dt$$

where $\Psi(\lambda) = \log \mathbb{E}[e^{\lambda(Y-\mathbb{E}Y)}]$. Then by Hoeffding's lemma, $\Psi''(\lambda) \leq \frac{(b-a)^2}{4}$. Thus

$$\int_0^\lambda t \Psi''(t) dt \leq \int_0^\lambda t \frac{(b-a)^2}{4} dt = \frac{(b-a)^2 \lambda^2}{8}.$$

\square

Log-Sobolev Inequalities

We want an analogue to the Poincaré inequality for entropy. Let X_1, \dots, X_n be independent symmetric Bernoulli random variables and let $f : \{-1, 1\}^n \rightarrow \mathbb{R}$. Then by Efron-Stein

$$\text{Var}(f(X)) \leq \frac{1}{2} \sum_{i=1}^n \mathbb{E}[(Z - Z'_i)^2]$$

where $Z = f(X)$ and $Z'_i = f(X^{(i)}, X'_i)$ for X'_i an iid copy of X_i , independent of all of X_1, \dots, X_n . Hence

$$\begin{aligned} \text{Var}(f(X)) &\leq \frac{1}{4} \sum_{i=1}^n \mathbb{E}[(f(X) - f(\bar{X}^{(i)}))^2] \\ &= \frac{1}{2} \sum_{i=1}^n \mathbb{E}[(f(X) - f(\bar{X}^{(i)}))_+^2] \end{aligned}$$

where $\bar{X}^{(i)} = f(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n)$. So define

$$\frac{1}{4} \sum_{i=1}^n \mathbb{E}[(f(X) - f(\bar{X}^{(i)}))^2] =: \mathcal{E}(f).$$

We think of $\mathcal{E}(f)$ as the expectation of a discrete derivative.

Theorem (Log-Sobolev inequality for symmetric Bernoulli).

$$\text{Ent}(f(X)^2) \leq 2\mathcal{E}(f).$$

Proof. Using tensorisation of Ent we have

$$\text{Ent}(Z^2) \leq \mathbb{E} \left[\sum_{i=1}^n \text{Ent}^{(i)}(Z^2) \right]$$

where $\text{Ent}^{(i)}(Z^2) = \mathbb{E}^{(i)}[Z^2 \log Z^2] - \mathbb{E}^{(i)}[Z^2] \log \mathbb{E}^{(i)}[Z^2]$. So if the inequality is true for $n = 1$, we have

$$\text{Ent}^{(i)}(Z^2) \leq \frac{(f(X) - f(\bar{X}^{(i)}))^2}{2}$$

almost-surely, so by summing over $1 \leq i \leq n$ and taking expectations we get the general result.

Now to show the inequality for $n = 1$, we will need to show it holds for $f(-1) = a, f(1) = b$. Indeed in this case

$$\text{Ent}(Z^2) = \frac{1}{2}a^2 \log a^2 + \frac{1}{2}b^2 \log b^2 - \frac{a^2 + b^2}{2} \log \left(\frac{a^2 + b^2}{2} \right)$$

and $2\mathcal{E}(f) = \frac{(b-a)^2}{2}$. So we need to show

$$\frac{1}{2}a^2 \log a^2 + \frac{1}{2}b^2 \log b^2 - \frac{a^2 + b^2}{2} \log \left(\frac{a^2 + b^2}{2} \right) \leq \frac{(b-a)^2}{2}.$$

WLOG $0 \leq b \leq a$. For fixed b consider $h : [b, \infty) \rightarrow \mathbb{R}$ defined by

$$h(a) = \frac{1}{2}a^2 \log a^2 + \frac{1}{2}b^2 \log b^2 - \frac{a^2 + b^2}{2} \log \left(\frac{a^2 + b^2}{2} \right) - \frac{(b-a)^2}{2}.$$

We have

$$\begin{aligned} h'(a) &= a \log \frac{2a^2}{a^2 + b^2} - (a - b) \\ h''(a) &= 1 + \log \frac{2a^2}{a^2 + b^2} - \frac{2a^2}{a^2 + b^2} \leq 0 \quad (\log x - x \leq -1) \end{aligned}$$

so $h(b) = h'(b) = 0$ and $h''(a) \leq 0$ for all $a \in [b, \infty)$, which implies $h(a) \leq 0$ for all $a \in [b, \infty)$. \square

Theorem (Log-Sobolev inequality for Gaussians). *Let X_1, \dots, X_n be iid $\mathcal{N}(0, 1)$, let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function. Then*

$$\text{Ent}(f^2) \leq 2\mathbb{E}[\|\nabla f(X)\|^2].$$

Proof sketch. Go through the following steps

1. Reduce it to the $n = 1$ case by tensorisation;
2. Introduce X_1, \dots, X_n iid symmetric Bernoulli's and consider $f\left(\frac{X_1 + \dots + X_n}{\sqrt{n}}\right)$ and use the log-Sobolev inequality for symmetric Bernoulli's;
3. Take $n \rightarrow \infty$ and use the CLT.

[Details of proof on Example Sheet] \square

Theorem (Gaussian concentration inequality). *Let X_1, \dots, X_n be iid $\mathcal{N}(0, 1)$. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be an L -Lipschitz function. Then $Z = f(X_{1:n})$ is sub-Gaussian with variance parameter L^2 .*

Proof. Apply the Gaussian log-Sobolev inequality to $e^{\lambda Z/2}$ to get

$$\text{Ent}(e^{\lambda Z}) \leq 2\mathbb{E}[\|e^{\lambda Z/2} \frac{\lambda}{2} \nabla f(X)\|^2] \leq \frac{\lambda^2}{2} L^2 \mathbb{E}[e^{\lambda Z}]$$

implying

$$\frac{\text{Ent}(e^{\lambda Z})}{\mathbb{E}(e^{\lambda Z})} \leq \frac{\lambda^2 L^2}{2} \implies Z \in \mathcal{G}(L^2).$$

\square

Theorem. Let $f : \{-1, 1\}^n \rightarrow \mathbb{R}$ and let X_i be iid symmetric Bernoulli. Let $Z = f(X_{1:n})$ and let

$$v = \max_{x \in \{-1, 1\}^n} \sum_{i=1}^n (f(x) - f(\bar{x}^{(i)}))_+^2.$$

Then Z has a sub-Gaussian right tail with parameter $v/2$, i.e

$$\mathbb{P}(Z - \mathbb{E}Z \geq t) \leq e^{-\frac{t^2}{2}}.$$

Remarks.

1. $\text{Var}(Z) \leq \mathcal{E}(f) \leq \frac{v}{2}$;
2. If $v = \max_{x \in \{-1, 1\}^n} \sum_{i=1}^n (f(x) - f(\bar{x}^{(i)}))_-^2$, get left tail bounds that are $\mathcal{G}(v/2)$;
3. If $v = \max_{x \in \{-1, 1\}^n} \sum_{i=1}^n (f(x) - f(\bar{x}^{(i)}))^2$, get right and left tail which are $\mathcal{G}(v/2)$. More refined analysis shows it can actually be made $\mathcal{G}(v/4)$;
4. If f satisfies bounded differences with c_i such that $\sum_{i=1}^n c_i^2 \leq v$ then bounded differences gives $Z \in \mathcal{G}(v/4)$. The bound above also gives $Z \in \mathcal{G}(v/4)$, but is applicable more generally.

Proof. Let $\lambda > 0$. Use log-Sobolev inequality for $e^{\lambda Z/2}$ to get

$$\text{Ent}(e^{\lambda Z}) \leq \mathbb{E} \left[\sum_{i=1}^n \left(e^{\lambda f(X)/2} - e^{\lambda f(\bar{X}^{(i)})/2} \right)_+^2 \right].$$

Since $x \mapsto e^{x/2}$ is convex and so if $z > y$ we have $e^{z/2} - e^{y/2} \leq (z - y) \frac{e^{z/2}}{2}$. Therefore

$$\begin{aligned} \mathbb{E} \left[\sum_{i=1}^n \left(e^{\lambda f(X)/2} - e^{\lambda f(\bar{X}^{(i)})/2} \right)_+^2 \right] &\leq \mathbb{E} \left[\sum_{i=1}^n \left(\lambda f(X) - \lambda f(\bar{X}^{(i)}) \right)_+^2 \frac{e^{\lambda f(X)}}{4} \right] \\ &\leq \mathbb{E} \left[\frac{e^{\lambda f(X)} \lambda^2}{4} v \right] \\ &= \mathbb{E}[e^{\lambda Z}] \frac{\lambda^2 (v/2)}{2}. \end{aligned}$$

Use Herbst's argument to get the right tail bound. \square

Theorem (Modified log-Sobolev inequality). Let X_1, \dots, X_n be independent, $f : \mathcal{X}^n \rightarrow \mathbb{R}$, $Z = f(X_{1:n})$. For $1 \leq i \leq n$ let $Z_i = f_i(X^{(i)})$. Let $\phi(x) = e^x - x - 1$. Then for all $\lambda \in \mathbb{R}$

$$\text{Ent}(e^{\lambda Z}) \leq \sum_{i=1}^n \mathbb{E} [e^{\lambda Z} \phi(-\lambda(Z - Z_i))].$$

Remark. If $x \geq 0$ then $\phi(-x) \leq \frac{x^2}{2}$. Say $\lambda > 0$, we choose Z_i so that $Z - Z_i \geq 0$. Then

$$\phi(-\lambda(Z - Z_i)) \leq \frac{\lambda^2}{2}(Z - Z_i)^2$$

so the RHS of the inequality becomes $\mathbb{E} \left[e^{\lambda Z} \frac{\lambda^2}{2} \sum_{i=1}^n (Z - Z_i)^2 \right]$.

Before we prove the theorem, we will need the following lemma.

Lemma (Variational formula for Ent). *Let $Y \geq 0$ almost-surely. Then $\text{Ent}(Y) = \inf_{u>0} \mathbb{E}[Y \log(Y/u) - (Y - u)]$.*

Remark. $\text{Var}(Y) = \inf_u \mathbb{E}[(Y - u)^2]$. In general

$$\mathbb{E}[\phi(Y)] - \Phi(\mathbb{E}[Y]) = \inf_u \underbrace{\mathbb{E}[\Phi(Y) - \Phi(u) - \Phi'(u)(Y - u)]}_{\text{Bregman divergence}}.$$

In particular taking $\phi(x) = x^2$ gives variance, and taking $\phi(x) = x \log x$ gives the lemma.

Proof. Taking $u = \mathbb{E}Y$ gives $\mathbb{E}[Y \log(Y/u) - (Y - u)] = \text{Ent}(Y)$. Suppose $\mathbb{E}Y = m$, fix some $u > 0$. We want to show

$$\begin{aligned} \mathbb{E}[Y \log(Y/u) - (Y - u)] &\geq \mathbb{E}[Y \log Y] - m \log m \\ \iff -m \log u - (m - u) &\geq -m \log m \\ \iff \log(m/u) &\geq 1 - \frac{u}{m} \end{aligned}$$

which is true since $-\log(x) \geq 1 - x$. □

Now we prove the theorem.

Proof of modified log-Sobolev inequality. Let $Y = e^{\lambda Z}$, $Y_i = e^{\lambda Z_i}$. Then

$$\begin{aligned} \text{Ent}(Y) &\leq \mathbb{E} \left[\sum_{i=1}^n \text{Ent}^{(i)}(Y) \right] \\ &\leq \mathbb{E} \left[\sum_{i=1}^n \mathbb{E}^{(i)} [e^{\lambda Z} \lambda (Z - Z_i) - (e^{\lambda Z} - e^{\lambda Z_i})] \right] \\ &= \sum_{i=1}^n \mathbb{E} [e^{\lambda Z} \phi(-\lambda(Z - Z_i))]. \end{aligned}$$

□

Theorem. Let $Z = f(X_{1:n})$ for independent X_1, \dots, X_n . Define $Z_i = \inf_{x_i} f(X^{(i)}, x_i)$. Suppose $\sum_{i=1}^n (Z - Z_i)^2 \leq v$. Then for all $t > 0$,

$$\mathbb{P}(Z - \mathbb{E}Z \geq t) \leq e^{-\frac{t^2}{2v}}.$$

Proof. By the modified log-Sovolev inequality

$$\begin{aligned} \text{Ent}(e^{\lambda Z}) &\leq \mathbb{E} \left[\sum_{i=1}^n e^{\lambda Z} \phi(-\lambda(Z - Z_i)) \right] \\ &\leq \mathbb{E} \left[\sum_{i=1}^n e^{\lambda Z} \frac{\lambda^2 (Z - Z_i)^2}{2} \right] \\ &\leq \frac{\lambda^2 v}{2} \mathbb{E}[e^{\lambda Z}]. \end{aligned}$$

So use Herbst's argument. \square

Theorem. Let f be a separately convex function on $[0, 1]^n$. Let X_1, \dots, X_n be independent and supported on $[0, 1]$. Let $Z = f(X_{1:n})$. Assume that f is 1-Lipschitz. Then $\mathbb{P}(Z - \mathbb{E}Z \geq t) \leq e^{-t^2/2}$ for $t > 0$.

Remark. $\text{Var}(Z) \leq 1$ by the convex Poincaré inequality.

Proof. Set $Z_i = \inf_{x'_i} f(X^{(i)}, x'_i)$. Let x_i^* be such that $Z_i = f(X^{(i)}, x_i^*)$. Then

$$\begin{aligned} Z_i &\geq Z + \frac{\partial f}{\partial x_i}(X) \cdot (x_i^* - X_i) \\ \implies 0 &\leq Z - Z_i \leq \frac{\partial f}{\partial x_i}(X) \cdot (X_i - x_i^*) \\ \implies (Z - Z_i)^2 &\leq \left(\frac{\partial f}{\partial x_i}(X) \right)^2. \end{aligned}$$

Summing up we get $\sum_{i=1}^n (Z - Z_i)^2 \leq \|\nabla f(X)\|^2 \leq 1$. Using the previous theorem we get $\mathbb{P}(Z - \mathbb{E}Z \geq t) \leq e^{-t^2/2}$. \square

Transport Method

We consider two sets of nodes $(x_i)_{i \in I}$ and $(y_j)_{j \in J}$ and we want to ‘transport’ materials from the (x_i) to the (y_j) . Each x_i has a production level and each y_j has a capacity. Furthermore each pair (x_i, y_j) has an associated transport cost $c(x_i, y_j)$. Our problem is to transport all the production from the x_i to the y_j , while minimising the total cost.

Definition. A *transport plan* is a function $\Pi(x_i, y_j)$ for $i \in I, j \in J$, where $\Pi(x_i, y_j)$ represents the amount transported from x_i to y_j . We define

$$\sum_y \pi(x_i, y) =: p(x_i)$$

$$\sum_x \pi(x, y_j) =: q(y_j)$$

and the optimal cost is

$$\min_{\Pi} \sum_{i,j} c(x_i, y_j) \Pi(x_i, y_j).$$

Theorem (Variational formulas for log-MGF and KL-divergence). *Let Z be a real valued random variable on a probability space (Ω, \mathcal{F}, P) . Then*

$$\log \mathbb{E}_P e^Z = \sup_{Q \ll P} [\mathbb{E}_Q Z - D(Q \| P)].$$

Conversely, if P and Q are two measures then

$$D(Q \| P) = \sup_Z [\mathbb{E}_Q Z - \log \mathbb{E}_P e^Z].$$

Remark. If Z is replaced by $\lambda(Z - \mathbb{E}_P Z)$ then

$$\log \mathbb{E}_P e^{\lambda(Z - \mathbb{E}_P Z)} = \sup_{Q \ll P} [\lambda(\mathbb{E}_Q Z - \mathbb{E}_P Z) - D(Q \| P)].$$

Proof. We assume Ω is discrete. Set $Q^*(\omega) = \frac{e^{Z(\omega)} P(\omega)}{\mathbb{E}_P e^Z}$. We have

$$\begin{aligned} 0 \leq D(Q \| Q^*) &\leq \sum_{\omega \in \Omega} Q(\omega) \log \frac{Q(\omega)}{Q^*(\omega)} \\ &= \sum_{\omega \in \Omega} Q(\omega) \log \left(\frac{Q(\omega)}{P(\omega)} \times \frac{P(\omega)}{Q^*(\omega)} \right) \\ &= D(Q \| P) + \sum_{\omega \in \Omega} Q(\omega) \log \frac{\mathbb{E}_P e^Z}{e^{Z(\omega)}} \\ &= D(Q \| P) + \log \mathbb{E}_P e^Z - \mathbb{E}_Q Z. \end{aligned}$$

Thus $\log \mathbb{E}_P e^Z \geq \mathbb{E}_Q Z - D(Q \| P)$. Furthermore Q^* achieves equality and $Q^* \ll P$, so we have shown the first part.

Conversely, by the above

$$D(Q \| P) \geq \mathbb{E}_Q Z - \log \mathbb{E}_P e^Z$$

and $Z(\omega) = \frac{Q(\omega)}{P(\omega)}$ gives equality. □

Suppose that we have a relationship of the form

$$\mathbb{E}_Q Z - \mathbb{E}_P Z \leq \sqrt{2vD(Q\|P)}$$

for all $Q \ll P$. Then by the above

$$\begin{aligned} \log \mathbb{E}_P e^{\lambda(Z - \mathbb{E}_P Z)} &= \sup_{Q \ll P} [\lambda(\mathbb{E}_Q Z - \mathbb{E}_P Z) - D(Q\|P)] \\ &\leq \sup_{Q \ll P} [\lambda\sqrt{2vD(Q\|P)} - D(Q\|P)] \\ &= \sup_{t \geq 0} \lambda\sqrt{2vt} - t \\ &= \frac{\lambda^2 v}{2}. \end{aligned}$$

Theorem (Marton's argument). *Suppose that for some $v > 0$, the following holds for all $Q \ll P$:*

$$\mathbb{E}_Q Z - \mathbb{E}_P Z \leq \sqrt{2vD(Q\|P)}.$$

Then for all $\lambda > 0$

$$\log \mathbb{E} e^{\lambda(Z - \mathbb{E}_P Z)} \leq \frac{\lambda^2 v}{2}$$

and $\mathbb{P}(Z - \mathbb{E} Z \geq t) \leq e^{-\frac{t^2}{2v}}$.

Conversely, if $\log \mathbb{E}_P e^{\lambda(Z - \mathbb{E}_P Z)} \leq \frac{\lambda^2 v}{2}$ for all $\lambda > 0$, then $\mathbb{E}_Q Z - \mathbb{E}_P Z \leq \sqrt{2vD(Q\|P)}$ for all $Q \ll P$.

Proof. We have

$$\begin{aligned} \log \mathbb{E}_P e^{\lambda(Z - \mathbb{E}_P Z)} &= \sup_{Q \ll P} [\lambda(\mathbb{E}_Q Z - \mathbb{E}_P Z) - D(Q\|P)] \\ &\leq \sup_{Q \ll P} [\lambda \sqrt{2vD(Q\|P)} - D(Q\|P)] \quad (\lambda > 0) \\ &\leq \sup_{t \geq 0} [\lambda \sqrt{2vt} - t] \\ &= \frac{\lambda^2 v}{2}. \end{aligned}$$

For the converse, assume $\mathbb{E}_Q Z - \mathbb{E}_P Z > 0$ (otherwise the result is trivial). We have (by taking $\lambda(Z - \mathbb{E}_P Z)$ in our variational formula for $D(Q\|P)$)

$$\begin{aligned} D(Q\|P) &\geq \lambda(\mathbb{E}_Q Z - \mathbb{E}_P Z) \log \mathbb{E}_P e^{\lambda(Z - \mathbb{E}_P Z)} \\ &\geq \lambda \lambda(\mathbb{E}_Q Z - \mathbb{E}_P Z) - \frac{\lambda^2 v}{2} \end{aligned}$$

so by maximising over λ , giving $\lambda = \frac{\mathbb{E}_Q Z - \mathbb{E}_P Z}{v}$ we get

$$D(Q\|P) \geq \frac{(\mathbb{E}_Q Z - \mathbb{E}_P Z)^2}{2v} \implies \mathbb{E}_Q Z - \mathbb{E}_P Z \leq \sqrt{2vD(Q\|P)}.$$

□

Suppose $X_{1:n} \sim P = P_{X_1} \otimes \dots \otimes P_{X_n}$, $Z = f(X_{1:n})$. Suppose further

$$f(y) - f(x) \leq \sum_{i=1}^n d(x_i, y_i) c_i$$

for some metric d (note that $d(x_i, y_i) := \mathbb{1}\{x_i \neq y_i\}$ gives bounded differences). Let $Y_{1:n} \sim Q$ (not necessarily a product distribution). Then

$$\mathbb{E} f(Y_{1:n}) - \mathbb{E} f(X_{1:n}) = \mathbb{E}_\Pi [f(Y_{1:n}) - f(X_{1:n})]$$

where the expectation on the RHS is taken with respect to a *coupling measure* Π between $X_{1:n}$ and $Y_{1:n}$, such that Π has $X_{1:n}$ -marginal P and $Y_{1:n}$ -marginal Q . Then

$$\begin{aligned}\mathbb{E}_\Pi[f(Y_{1:n}) - f(X_{1:n})] &\leq \mathbb{E}_\Pi \left[\sum_{i=1}^n d(X_i, Y_i) c_i \right] \\ &= \sum_{i=1}^n c_i \mathbb{E}[d(X_i, Y_i)] \\ &\leq \left(\sum_{i=1}^n c_i^2 \right)^{1/2} \left(\sum_{i=1}^n [\mathbb{E}_\Pi[d(X_i, Y_i)]^2] \right)^{1/2}\end{aligned}$$

now note we can actually optimise over couplings Π to get

$$\mathbb{E}[f(Y_{1:n})] - \mathbb{E}[f(X_{1:n})] \leq \left(\sum_{i=1}^n c_i^2 \right)^{1/2} \left(\inf_{\Pi \in \Pi(P, Q)} \sum_{i=1}^n [\mathbb{E}_\Pi[d(X_i, Y_i)]^2] \right)^{1/2}$$

where $\Pi(P, Q)$ denotes the set of all couplings between P and Q . Suppose we show

$$\inf_{\Pi \in \Pi(P, Q)} \sum_{i=1}^n [\mathbb{E}_\Pi[d(X_i, Y_i)]^2] \leq 2CD(Q\|P)$$

then we would have

$$\mathbb{E}[f(Y_{1:n})] - \mathbb{E}[f(X_{1:n})] \leq \sqrt{2vD(Q\|P)}$$

where $v = C \sum_{i=1}^n c_i^2$.

Therefore, to run Marton's argument for such functions f it is enough to prove

$$\inf_{\Pi} \sum_{i=1}^n \mathbb{E}_\Pi[d(X_i, Y_i)]^2 \leq 2CD(Q\|P) \text{ for some } C > 0.$$

Bounded differences inequality via the transport method

We want to show

$$\begin{aligned}\inf_{\Pi} \sum_{i=1}^n \mathbb{E}_\Pi[\mathbb{1}(X_i \neq Y_u)]^2 &= \inf_{\Pi} \sum_{i=1}^n \mathbb{P}_\Pi(X_i \neq Y_u)^2 \\ &\leq \frac{1}{2} D(Q\|P)\end{aligned}$$

which will then give us the bounded differences inequality we had before.

Theorem (Marton's transport cost inequality). *Let $P = P_{X_1} \otimes \dots \otimes P_{X_n}$ and Q be an arbitrary measure over the same measurable space such that $Q \ll P$. Then*

$$\inf_{\Pi \in \Pi(P, Q)} \sum_{i=1}^n \mathbb{P}_\Pi(X_i \neq Y_u)^2 \leq \frac{1}{2} D(Q\|P).$$

Remark. For $n = 1$ we would have $\inf_{\Pi} \mathbb{P}(X_1 \neq Y_1) \leq \sqrt{D(Q\|P)/2}$.

Lemma. Let P, Q be probability measures on the same measurable space (Ω, \mathcal{F}) . Then

$$\inf_{(X,Y) \sim \Pi \in \Pi(P,Q)} \mathbb{P}(X \neq Y) = d_{TV}(P, Q)$$

where $d_{TV}(P, Q)$ is the total variation distance between P and Q , defined by

$$d_{TV}(P, Q) = \sup_{A \in \mathcal{F}} |P(A) - Q(A)|.$$

Remark. We can also write

$$\begin{aligned} d_{TV}(P, Q) &= \frac{1}{2} \sum_{\omega \in \Omega} |P(\omega) - Q(\omega)| \\ &= \sum_{\omega \in \Omega} (P(\omega) - Q(\omega))_+ \\ &= 1 - \sum_{\omega \in \Omega} \min\{P(\omega), Q(\omega)\}. \end{aligned}$$

Proof. For any $A \in \mathcal{F}$ we have

$$\begin{aligned} |P(A) - Q(A)| &= |\mathbb{E}_{\Pi}[\mathbb{1}(X \in A)] - \mathbb{E}_{\Pi}[\mathbb{1}(Y \in A)]| \\ &\leq \mathbb{E}_{\Pi}[\mathbb{1}(X \neq Y)] \\ &= \mathbb{P}_{\Pi}(X \neq Y) \end{aligned}$$

so taking suprema over $A \in \mathcal{F}$ and infima over Π we get

$$d_{TV}(P, Q) \leq \inf_{\Pi} \mathbb{P}(X \neq Y).$$

So it remains to find $\Pi \in \Pi(P, Q)$ such that $\mathbb{P}(X \neq Y) = d_{TV}(P, Q)$. Let $A = \{\omega : P(\omega) \geq Q(\omega)\}$. Consider the coupling

$$\Pi(\omega_1, \omega_2) = \begin{cases} Q(\omega) & \omega_1 = \omega_2 \in A \\ P(\omega) & \omega_1 = \omega_2 \in A^c \\ 0 & (\omega_1, \omega_2) \in A^c \times A \\ \frac{(P(\omega_1) - Q(\omega_1))(Q(\omega_2) - P(\omega_2))}{d_{TV}(P, Q)} & (\omega_1, \omega_2) \in A \times A^c \end{cases}.$$

Then we have $\mathbb{P}_{\Pi}(X = Y) = \sum_{\omega} \min\{P(\omega), Q(\omega)\} = 1 - d_{TV}(P, Q)$ which implies $\mathbb{P}(X \neq Y) = d_{TV}(P, Q)$. \square

Lemma (Pinsker's inequality). *We have*

$$d_{TV}(P, Q)^2 \leq \frac{1}{2} D(Q \| P).$$

Proof. Example Sheet 2. □

Proof of Marton's transport cost inequality. The above lemmas prove Marton's TCI for $n = 1$. Assume that it holds for all $n \leq k$ and we prove it for $n = k + 1$. Let $(X_1, \dots, X_{k+1}) \sim P_{X_1} \otimes \dots \otimes P_{X_{k+1}}$ and $(Y_1, \dots, Y_{k+1}) \sim Q_{Y_{1:k+1}}$. We want to show

$$\inf_{\Pi \in \Pi(P_{X_{1:k+1}}, Q_{Y_{1:k+1}})} \sum_{i=1}^{k+1} \mathbb{P}(X_i \neq Y_i)^2 \leq \frac{1}{2} D(Q_{Y_{1:k+1}} \| P_{X_{1:k+1}}).$$

We know there exists $\Pi_k \in \Pi(P_{X_{1:k}}, Q_{Y_{1:k}})$ such that

$$\sum_{i=1}^k \mathbb{P}(X_i \neq Y_i)^2 \leq \frac{1}{2} D(Q_{Y_{1:k}} \| P_{X_{1:k}}).$$

Define $\Pi \in \Pi(P_{X_{1:k+1}}, Q_{Y_{1:k+1}})$ as

$$\begin{aligned} \Pi(X_{1:k+1} = x_{1:k+1}, Y_{1:k+1} = y_{1:k+1}) \\ = \Pi_k(X_{1:k} = x_{1:k}, Y_{1:k} = y_{1:k}) \Pi_{y_{1:k}}(X_{k+1} = x_{k+1}, Y_{k+1} = y_{k+1}) \end{aligned}$$

where $\Pi_{y_{1:k}}$ is the optimal TV-coupling between $P_{X_{k+1}}$ and $Q_{Y_{k+1}|Y_{1:k}=y_{1:k}}$. Under Π ,

$$\begin{aligned} \mathbb{P}(X_{1:k+1} = x_{1:k+1}, Y_{1:k+1} = y_{1:k+1}) \\ = \mathbb{P}(X_{1:k} = x_{1:k}, Y_{1:k} = y_{1:k}) \times \mathbb{P}(X_{k+1} = x_{k+1}) \mathbb{P}(Y_{k+1} = y_{k+1} | Y_{1:k} = y_{1:k}, X_{k+1} = x_{k+1}). \end{aligned}$$

Under Π we have

$$\sum_{i=1}^k \mathbb{P}(X_i \neq Y_i)^2 + \mathbb{P}(X_{k+1} \neq Y_{k+1})^2 \leq \frac{1}{2} D(Q_{Y_{1:k}} \| P_{X_{1:k}}) + \mathbb{P}(X_{k+1} \neq Y_{k+1})^2.$$

Observe that

$$\begin{aligned} \mathbb{P}_{\Pi}(X_{k+1} \neq Y_{k+1} | X_{1:k} = x_{1:k}, Y_{1:k} = y_{1:k}) &= \mathbb{P}_{\Pi}(X_{k+1} \neq Y_{k+1} | Y_{1:k} = y_{1:k}) \\ &= d_{TV}(P_{X_{k+1}}, Q_{Y_{k+1}|Y_{1:k}=y_{1:k}}) \\ &\leq \sqrt{\frac{1}{2} D(Q_{Y_{k+1}|Y_{1:k}=y_{1:k}} \| P_{X_{k+1}})} \end{aligned}$$

so integrate with respect to Π_k to see

$$\mathbb{P}(X_{k+1} \neq Y_{k+1}) \leq \mathbb{E}_{\Pi_k} \left[\sqrt{\frac{1}{2} D(Q_{Y_{k+1}|Y_{1:k}=y_{1:k}} \| P_{X_{k+1}})} \right]$$

so squaring and applying Jensen,

$$\begin{aligned}\mathbb{P}(X_{k+1} \neq Y_{k+1}) &\leq \frac{1}{2} \mathbb{E}_{\Pi_k} D(Q_{Y_{k+1}|Y_{1:k}=y_{1:k}} \| P_{X_{k+1}}) \\ &= \frac{1}{2} \mathbb{E}_{Q_{Y_{1:k}}} D(Q_{Y_{k+1}|Y_{1:k}=y_{1:k}} \| P_{X_{k+1}}) \\ &= \frac{1}{2} D(Q_{Y_{k+1}|Y_{1:k}} \| P_{X_{k+1}} | Q_{Y_{1:k}}).\end{aligned}$$

Therefore

$$\begin{aligned}&\sum_{i=1}^k \mathbb{P}_{\Pi}(X_i \neq Y_i)^2 + \mathbb{P}(X_{k+1} \neq Y_{k+1})^2 \\ &\leq \frac{1}{2} D(Q_{Y_{1:k}} \| P_{X_{1:k}}) + \frac{1}{2} D(Q_{Y_{k+1}|Y_{1:k}} \| P_{X_{k+1}} | Q_{Y_{1:k}}) \\ &= \frac{1}{2} D(Q_{Y_{1:k+1}} \| P_{X_{1:k+1}}). \quad (\text{chain rule})\end{aligned}$$

□

Definition. A function $f : \mathcal{X}^n \rightarrow \mathbb{R}$ satisfies a *one-sided bounded differences property* with functions c_1, \dots, c_n from \mathcal{X}^n to \mathbb{R} if for all $x, y \in \mathcal{X}^n$ we have

$$f(y) - f(x) \leq \sum_{i=1}^n c_i(x) \mathbb{1}\{x_i \neq y_i\}.$$

Theorem (Talagrand's one-sided bounded differences inequality). *Let X_1, \dots, X_n be independent and let f satisfy one-sided bounded differences with $c_1, \dots, c_n : \mathcal{X}^n \rightarrow \mathbb{R}$. Define $v = \mathbb{E} [\sum_{i=1}^n c_i(X)^2]$ and $Z = f(X_{1:n})$. Then for $\lambda > 0$*

$$\psi_{Z - \mathbb{E}Z}(\lambda) \leq \frac{\lambda^2 v}{2}$$

and so for $t > 0$, $\mathbb{P}(Z - \mathbb{E}Z \geq t) \leq e^{-\frac{t^2}{2v}}$.

Remark. If $v_\infty := \sup_x \sum_{i=1}^n c_i(x)^2$ and $Z_i = \sup_{x_i} f(X^{(i)}, x_i)$, one-sided bounded differences gives

$$0 \leq Z_i - Z \leq c_i(X) \implies \sum_{i=1}^n (Z_i - Z)^2 \leq v_\infty$$

so we get a left tail with parameter v_∞ .

Proof. Let $P = P_{X_1} \otimes \dots \otimes P_{X_n}$ and let $Y_{1:n} \sim Q$. We have $f : \mathcal{X}^n \rightarrow \mathbb{R}$ and we want to bound

$$\mathbb{E}f(Y_{1:n}) - \mathbb{E}f(X_{1:n}) = \mathbb{E}_\pi[f(Y_{1:n}) - f(X_{1:n})]$$

for $\pi \in \Pi(P, Q)$. We have

$$\begin{aligned} \mathbb{E}_\pi[f(Y_{1:n}) - f(X_{1:n})] &\leq \mathbb{E}_\pi \left[\sum_{i=1}^n c_i(X_{1:n}) \mathbb{1}\{X_i \neq Y_i\} \right] \\ &= \mathbb{E}_\pi \left[\mathbb{E}_\pi \left[\sum_{i=1}^n c_i(X_{1:n}) \mathbb{1}\{X_i \neq Y_i\} \mid X_{1:n} \right] \right] \\ &= \mathbb{E}_\pi \left[\sum_{i=1}^n c_i(X_{1:n}) \mathbb{P}(X_i \neq Y_i \mid X_{1:n}) \right] \\ &\leq \mathbb{E}_\pi \left[\left(\sum_{i=1}^n c_i(X_{1:n})^2 \right)^{1/2} \left(\sum_{i=1}^n \mathbb{P}(X_i \neq Y_i \mid X_{1:n}) \right)^{1/2} \right] \\ &\leq \sqrt{\mathbb{E}_\pi \left[\sum_{i=1}^n c_i(X_{1:n})^2 \right]} \sqrt{\mathbb{E}_\pi \left[\sum_{i=1}^n \mathbb{P}(X_i \neq Y_i \mid X_{1:n}) \right]} \\ &= \sqrt{v} \left(\sum_{i=1}^n \mathbb{E}_\pi [\mathbb{P}(X_i \neq Y_i \mid X_{1:n})^2] \right)^{1/2}. \end{aligned}$$

Therefore it is enough to show that

$$\inf_{\pi \in \Pi(P, Q)} \sum_{i=1}^n \mathbb{E}_\pi [\mathbb{P}(X_i \neq Y_i \mid X_{1:n})^2] \leq 2D(Q \| P).$$

Claim (Marton's conditional transport cost inequality). *We have*

$$\inf_{\pi \in \Pi(P, Q)} \sum_{i=1}^n \mathbb{E}_\pi [\mathbb{P}(X_i \neq Y_i \mid X_{1:n})^2] \leq 2D(Q \| P).$$

We first show the $n = 1$ case:

Lemma. *Let P, Q be probability measures on a measurable space. Then*

$$\inf_{\pi \in \Pi(P, Q)} \mathbb{E}[\mathbb{P}(X \neq Y \mid X)^2] = d_2^2(Q, P)$$

where $d_2^2(Q, P)$ is Marton's divergence and is defined by

$$d_2^2(Q, P) = \sum_{\omega: P(\omega) > 0} \frac{(P(\omega) - Q(\omega))_+^2}{P(\omega)}.$$

Proof. Let π be any coupling. Observe that

$$\mathbb{P}(X = Y | X = x) = \frac{\mathbb{P}(X = x, Y = x)}{\mathbb{P}(X = x)} \leq \frac{\mathbb{P}(Y = x)}{\mathbb{P}(X = x)} = \frac{Q(x)}{P(x)}$$

so

$$\mathbb{P}(X \neq Y | X = x) \geq \left(1 - \frac{Q(x)}{P(x)}\right)_+.$$

Squaring this and taking expectations

$$\mathbb{E}_\pi[\mathbb{P}(X \neq Y | X)^2] \geq \sum_x P(x) \frac{(P(x) - Q(x))_+^2}{P(x)^2} = d_2^2(Q, P).$$

The existence of a coupling giving equality is on the example sheet. \square

Lemma. *We have*

$$d_2^2(Q, P) \leq 2D(Q \| P)$$

Proof. Omitted and non-examinable. \square

The above lemma implies Marton's conditional transport cost inequality for $n = 1$. We will use induction. Assume the result is true for $n \leq k$. Then we want to show

$$\inf_\pi \mathbb{E}_\pi[\text{left}[\sum_{i=1}^{k+1} \mathbb{P}(X_i \neq Y_i | X_{1:k+1})^2]] \leq D(Q_{Y_{1:k+1}} \| P_{X_{1:k+1}}).$$

We know there exists a coupling $\pi_k \in \Pi(P_{X_{1:k}}, Q_{Y_{1:k}})$ such that

$$\mathbb{E}_{\pi_k} \left[\sum_{i=1}^k \mathbb{P}(X_i \neq Y_i | X_{1:k})^2 \right] \leq 2D(Q_{Y_{1:k}} \| P_{X_{1:k}}).$$

Define

$$\begin{aligned} \pi(X_{1:k+1} = x_{1:k+1}, Y_{1:k+1}) \\ = \pi_k(X_{1:k} = x_{1:k}, Y_{1:k} = y_{1:k}) \pi_{y_{1:k}}(X_{k+1} = x_{k+1}, Y_{k+1} = y_{k+1}) \end{aligned}$$

where $\pi_{y_{1:k}}$ is the optimal TV coupling between $P_{X_{k+1}}$ and $Q_{Y_{k+1} | Y_{1:k} = y_{1:k}}$. Then π has the following nice properties:

- π has marginal π_k on $(X_{1:k}, Y_{1:k})$;
- (X_{k+1}, Y_{k+1}) depend on $(X_{1:k}, Y_{1:k})$ only through $Y_{1:k}$;
- X_{k+1} is independent of $(X_{1:k}, Y_{1:k})$.

Note that $\mathbb{P}(X_i \neq Y_i | X_{1:k+1}) = \mathbb{P}(X_i \neq Y_i | X_{1:k})$ for all $1 \leq i \leq k$. By the assumption for $n \leq k$ we conclude

$$\mathbb{E}_\pi \left[\sum_{i=1}^k \mathbb{P}(X_i \neq Y_i | X_{1:k+1})^2 \right] \leq 2D(Q_{Y_{1:k}} \| P_{1:k}).$$

We know by the choice of $\pi_{y_{1:k}}$ that

$$\mathbb{E}_{\pi_{y_{1:k}}} [\mathbb{P}(X_{k+1} \neq Y_{k+1} | Y_{1:k} = y_{1:k}, X_{k+1})] \leq 2D(Q_{Y_{k+1} | Y_{1:k} = y_{1:k}} \| P_{X_{k+1}}).$$

by the $n = 1$ result. If we then integrate both sides of this inequality with respect to the $Q_{Y_{1:k}}$ measure we get

$$\mathbb{E}_\pi [\mathbb{P}(X_{k+1} \neq Y_{k+1} | Y_{1:k}, X_{k+1})^2] \leq 2D(Q_{Y_{k+1} | Y_{1:k}} \| P_{X_{k+1}} | Q_{Y_{1:k}}).$$

Observe that since the distribution of (X_{k+1}, Y_{k+1}) depends only on $Y_{1:k}$ given $(X_{1:k}, Y_{1:k})$

$$\mathbb{P}(X_{k+1} \neq Y_{k+1} | Y_{1:k}, X_{k+1}) = \mathbb{P}(X_{k+1} \neq Y_{k+1} | Y_{1:k}, X_{1:k+1})$$

and so

$$\begin{aligned} \mathbb{E}[\mathbb{P}(X_{k+1} \neq Y_{k+1} | Y_{1:k}, X_{1:k+1})^2] &= \mathbb{E}[\mathbb{E}[\mathbb{1}\{X_{k+1} \neq Y_{k+1}\} | X_{1:k+1}, Y_{1:k}]^2] \\ &= \mathbb{E}[\mathbb{E}[\mathbb{E}[\mathbb{1}\{X_{k+1} \neq Y_{k+1}\} | X_{1:k+1}, Y_{1:k}]^2 | X_{1:k+1}]] \\ &\geq \mathbb{E}[\mathbb{E}[\mathbb{E}[\mathbb{1}\{X_{k+1} \neq Y_{k+1}\} | X_{1:k+1}, Y_{1:k}] | X_{1:k+1}]^2] \\ &\quad \text{(Jensen)} \\ &= \mathbb{E}[\mathbb{E}[\mathbb{1}\{X_{k+1} \neq Y_{k+1}\} | X_{1:k+1}]^2] \\ &\quad \text{(Tower property)} \\ &= \mathbb{E}[\mathbb{P}(X_{k+1} \neq Y_{k+1} | X_{1:k+1})^2] \end{aligned}$$

so putting this together

$$\begin{aligned} \mathbb{E}[\mathbb{P}(X_{k+1} \neq Y_{k+1} | X_{1:k+1})^2] &\leq \mathbb{E}[\mathbb{P}(X_{k+1} \neq Y_{k+1} | Y_{1:k}, X_{1:k+1})^2] \\ &\leq 2D(Q_{Y_{k+1} | Y_{1:k}} \| P_{X_{k+1}} | Q_{Y_{1:k}}). \end{aligned}$$

And thus

$$\begin{aligned} &\mathbb{E}_\pi \left[\sum_{i=1}^k \mathbb{P}(X_i \neq Y_i | X_{1:k+1})^2 \right] + \mathbb{E}[\mathbb{P}(X_{k+1} \neq Y_{k+1} | X_{1:k+1})^2] \\ &\leq 2D(Q_{Y_{1:k}} \| P_{1:k}) + 2D(Q_{Y_{k+1} | Y_{1:k}} \| P_{X_{k+1}} | Q_{Y_{1:k}}). \end{aligned}$$

□