[A] Minmax Redundancy and Probability Assignment

$\underline{x} \in \mathcal{Z}^n \rightsquigarrow (X_1, \ldots, X_n) = (x_1, \ldots, x_n)$ } probabilistic modelling

$X_i \sim P$, $(X_1, \ldots, X_n)$ are iid with common distribution P.

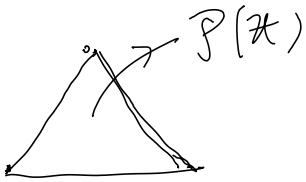→ Universal compression algorithms assume the generating dist. is unknown

→ Benchmark of performance: "Redundancy"

$\mathcal{C}$ be a prefix-free code that assigns a codeword of length $l(\underline{x})$

$$r_n(\mathcal{C}, P^n) := \sum_{\underline{x} \in \mathcal{Z}^n} P^n(\underline{x}) \, l(\underline{x}) - H(P^n) \longrightarrow$$

"redundancy of $\mathcal{C}$ for $P^n$"

Worst-case redundancy:


$\hat{P}(\mathcal{Z})$

$$r_n(\mathcal{C}) := \max_{P \in \mathcal{P}(\mathcal{Z})} r_n(\mathcal{C}, P^n)$$

→ This is our measure of for the prefix-free cod

Minmax redundancy

$$r_n^* := \min_{\mathcal{C}} r_n(\mathcal{C})$$

$$= \min_{\mathcal{C}} \max_{P \in \mathcal{P}(\mathcal{Z})} \left[ \sum_{\underline{x}} P^n(\underline{x}) \, l(\underline{x} \right.$$
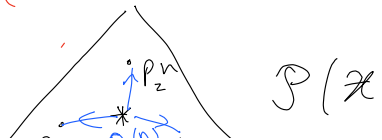
How large is $r_n^*$? And which codes $\mathcal{C}$ attain it?

Probability assignment $\rightarrow Q^{(n)}(\underline{x})$

We want distribution $Q^{(n)}$ on $\mathcal{Z}^n$ s.t.

$$R_n^* = \min_{Q^{(n)} \in \mathcal{P}(\mathcal{Z}^n)} \max_{P \in \mathcal{P}(\mathcal{Z})} D(P^n \| Q^{(n)})$$

(information rad


$P_z^n$
$\mathcal{P}(\mathcal{Z}$

Lemma. $R_n^* \approx r_n^*$

$$R_n^* \leq \eta_n^* \leq R_n^* + 1$$

**Proof.** (Probability assignment $\Rightarrow$ prefix-free code)

Given a prob. $Q^{(n)}$ on $\mathcal{X}^n$, let

$$\ell(\underline{x}) = \left\lceil \log \frac{1}{Q^{(n)}(\underline{x})} \right\rceil.$$

(check: $\ell(\underline{x})$ satisfies Kraft's inequality) ✓

Therefore, for the prefix-free code associated with $(\ell(\underline{x}), \underline{x}$
we have

$$\mathbb{E}_p[\ell(\underline{X})] = \sum_{\underline{x}} P^n(\underline{x}) \ell(\underline{x})$$

$$\leq \sum_{\underline{x}} P^n(\underline{x}) \log \frac{1}{Q^{(n)}(\underline{x})} \cdot \frac{P^n}{P^n}$$

$$= D(P^n \| Q^{(n)}) + nH(P) + 1$$

$$\Rightarrow \mathbb{E}_p[\ell(\underline{X})] - nH(P) \leq D(P^n \| Q^{(n)}) + 1$$

$$\Rightarrow \max_{P \in \mathcal{P}(\mathcal{X})} \mathbb{E}_p[\ell(\underline{X})] - nH(P) \leq \max_{P \in \mathcal{P}(\mathcal{X})} D(P^n \|$$

$$\Rightarrow \boxed{\eta_n^* \leq R_n^* + 1}$$

(prefix-free codes $\Rightarrow$ probability assignment)

Given a prefix-free code with codeword lengths $(\ell(\underline{x}), \underline{x} \in \mathcal{X}$
consider

$$Q^{(n)}(\underline{x}) = \frac{2^{-\ell(\underline{x})}}{\sum_{\underline{x}'} 2^{-\ell(\underline{x}')}}$$

Then,

$$D(P^n \| Q^{(n)}) = \sum_{\underline{x}} P^n(\underline{x}) \log \frac{1}{Q^{(n)}(\underline{x})} - nH(P)$$

$$= \sum_{\underline{x}} P^n(\underline{x}) \log 2^{\ell(\underline{x})} + \log \sum_{\underline{x}'} 2$$

$$\leq \mathbb{E}_{P^n}[\ell(\underline{X})] - n H(P),$$

for every $P \in \mathcal{P}(\mathcal{X})$. Therefore,

$$\min_{Q^{(n)} \in \mathcal{P}(\mathcal{X}^n)} \max_{P \in \mathcal{P}(\mathcal{X})} D\left(P^n \| Q^{(n)}\right) \leq \max_{P \in \mathcal{P}(\mathcal{X})} \mathbb{E}_{P^n}[$$

for every prefix-free code.
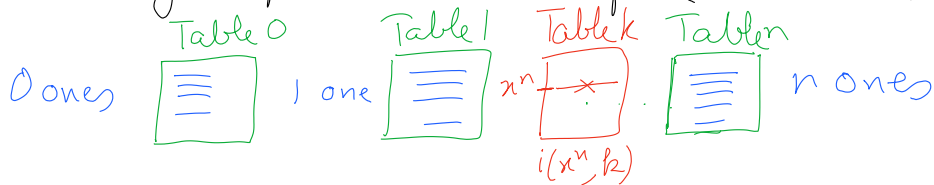
$$\Rightarrow \boxed{R_n^* \leq \mathcal{R}_n^*}$$

Universal prefix-free code design $\approx$ min max probability as

## [B] Compression using word frequencies

$\mathcal{X} = \{0, 1\}$

Our scheme :  Input: $x^n \in \{0, 1\}^n$

Output: A binary codeword $c(\underline{x}) \in \{0, 1\}^*$

1) Count the # of 1s in $x^n$. Denote it by $k$.

2) Let $i(x^n, k)$ be the index of the sequence $x^n$ among all sequen $k$ 1s.

(3) $c(\underline{x}) \equiv$ binary representation of $(k, i(x^n, k))$



What is the number bits used to represent a sequence $x^n$?

$\rightarrow$ $k$ takes $\log(n+1)$ bits to represent

$\rightarrow$ $i(x^n, k)$ can take $\binom{n}{k}$ values, and so, needs $\log\binom{n}{k}$ bi

Therefore,

$$\ell(\underline{x}) = \log(n+1) + \log\binom{n}{k(\underline{x})},$$

where $k(\underline{x}) = $ # 1s in $\underline{x}$.

**Important fact.** $\binom{n}{k} \approx 2^{n h(k/n)}$,

where $h(t) = t \log \frac{1}{t} + (1-t) \log \frac{1}{1-t}$.

More formally, $\binom{n}{k} \le 2^{n h\left(\frac{k}{n}\right)}$

$\Rightarrow \quad \ell(\underline{x}) \le \log(n+1) + n h\left(\frac{k(\underline{x})}{n}\right)$

Thus,

$$\mathbb{E}_{P^n}\left[\ell(X^n)\right] \le \log(n+1) + n \mathbb{E}_{P^n}\left[h\left(\frac{k(X}{r}\right)\right.$$

(by Jensen's ineq., since $h$ is a concave function) $\le \log(n+1) + n h\left(\frac{\mathbb{E}_{P^n}\left[k(X^n}{n}\right.\right.$

$$\mathbb{E}_{P^n}\left[k(X^n)\right] = \mathbb{E}_{P^n}\left[\sum_{t=1}^{n} X_t\right] = n p$$

$\Rightarrow \quad \mathbb{E}_{P^n}\left[\ell(X^n)\right] \le \log(n+1) + n \underbrace{h(p)}_{\substack{\longrightarrow P(1) \\ = H(P)}}$

$\Rightarrow \quad \mathbb{E}_{P^n}\left[\ell(X^n)\right] - n H(P) \le \log(n+1)$

for every $P \in \mathcal{P}(\{0,1\})$

$\Rightarrow \quad \max_{P \in \mathcal{P}(\mathcal{X}^n)} \mathbb{E}_{P^n}\left[\ell(X^n)\right] - n H(P) \le \log(n$

Remarks 1) $\log(n+1)$ extra cost for "universality" is negligible in comparison with the optimal avg. length $n H(P)$.

(2) Recall that we can find prob. assignment using this schem

$$R_n^* \le \mathfrak{r}_n^* \le \log(n+1)$$

$\longleftarrow$ iid $\Theta^{(n)}$ can only give $O(n)$ bounds

(3) The analysis above can be improved to get $\frac{1}{2} \log(n+1)$.

<span style="color:red">Extension to an arbitrary alphabet $\mathcal{X}$: Types (Method of Types, Csiszár - Körner b</span>

<u>Definition</u> (Type of a sequence) The <u>type of a sequence</u> $\underline{x}$ is a pmf denoted $P_{\underline{x}}$, given by

$$P_{\underline{x}}(a) = \frac{N(a|\underline{x})}{n}, \qquad a \in \mathcal{X}.$$

<span style="color:blue">$\rightarrow$ # of times $a$ appea $(x_1, \ldots, x_n)$</span>

The set of all sequences of a given type $Q$ is called the <u>type</u> type class $Q$, denoted $\mathcal{T}_Q^{(n)}$.

The set of all types is denoted by $T^{(n)}$.

<span style="color:blue"><u>Fact</u>: All sequences $\underline{x} \in \mathcal{T}_Q$ have equal probabilities under any iid</span>

<u>Proof</u>. 
$$P^n(\underline{x}) = \prod_{t=1}^{n} P(x_t)$$

$$= \prod_{a \in \mathcal{X}} P(a)^{N(a|\underline{x})}$$

$$= 2^{\sum_{a \in \mathcal{X}} N(a|\underline{x}) \log P(a)}$$

<span style="color:blue">$\rightarrow n Q(a)$</span>

$$= 2^{-\sum_{a \in \mathcal{X}} N(a|\underline{x}) \log \frac{1}{P(a)}} \cdot \frac{Q(a)}{Q(a)}$$

$$= 2^{-n\left(\sum_{a \in \mathcal{X}} Q(a) \log \frac{Q(a)}{P(a)} + H(Q)\right)}$$

$$= 2^{-n\left(D(Q||P) + H(Q)\right)}$$



<u>Lemma</u> (Type counting lemma) For a finite alphabet $\mathcal{X}$,

$$\boxed{|T^{(n)}| \leq (n+1)^{|\mathcal{X}|-1}}.$$

<u>Lemma</u> (Type class cardinality lemma) For every $Q \in T^{(n)}$,

$$\boxed{\frac{2^{nH(Q)}}{(n+1)^{|\mathcal{X}|-1}} \leq |\mathcal{T}_Q| \leq 2^{nH(Q)}}$$

$$\left( \frac{1}{n} \log |\mathcal{I}_Q| \approx H(Q) \right)$$

**Proof**

* $1 \geq Q^n(\mathcal{I}_Q) = \sum_{x \in \mathcal{I}_Q} Q^n(\underline{x}) = |\mathcal{I}_Q| \, 2^{-n H(Q)}$

$$\implies |\mathcal{I}_Q| \leq 2^{n H(Q)}.$$

* $\boxed{\displaystyle \max_{P \in T^{(n)}} Q^n(\mathcal{I}_P) = Q^n(\mathcal{I}_Q)}$  <span style="color:green">will show this</span>

<span style="color:green">Proof:</span>

$$\frac{Q^n(\mathcal{I}_P)}{Q^n(\mathcal{I}_Q)} = \frac{|\mathcal{I}_P|}{|\mathcal{I}_Q|} \cdot \frac{\prod_a Q(a)^{n P(a)}}{\prod_a Q(a)^{n Q(a)}}$$

$$= \frac{|\mathcal{I}_P|}{|\mathcal{I}_Q|} \cdot \prod_a Q(a)^{n(P(a) - Q(a))}$$

$$= \frac{n!}{\prod_a (n P(a))!} \cdot \frac{\prod_a (n Q(a))!}{n!} \quad \prod_a Q(a)^{n(P(a))}$$

$$= \frac{\prod_a (n Q(a))!}{(n P(a))!} \, Q(a)^{n(P(a) - Q(a))}$$

<span style="color:blue">Show this</span>

$$\boxed{\frac{k!}{\ell!} \leq k^{k-\ell}}$$

$$\leq \prod_a n Q(a)^{n(Q(a) - P(a) + P(a) - Q(a))}$$

$$= 1 \qquad \square$$

$$\to \quad 1 = \sum_{P \in T^{(n)}} Q^n(\mathcal{I}_P) \leq |T^{(n)}| \, Q^n(\mathcal{I}_Q)$$

$$\leq (n+1)^{|\mathcal{X}|-1} \cdot 2^{-n H(Q)} |\mathcal{I}_Q|$$

$$\not\Rightarrow \quad \boxed{|\mathcal{I}_Q| \geq \frac{2^{n H(Q)}}{(n+1)^{|\mathcal{X}|-1}}}$$

<span style="color:blue">General scheme  Given a sequence $\underline{x} \in \mathcal{X}^n$,

1) Find type $Q$ of $\underline{x}$

(2) Find the index $i(\underline{x}; Q)$ of $\underline{x}$ in the list of all sequences of $\mathcal{H}$</span>

[3] Store binary representation of $(Q, i(\underline{x}; Q))$

$\searrow$ $\quad$ $\log 2^r$

$\leq \log(n+1)^{|\mathcal{X}|-1}$ bits

$\Rightarrow l(\underline{x}) \leq (|\mathcal{X}|-1)\log(n+1) + n H(Q)$

Therefore, $\quad \mathbb{E}_{p^n}\left[l(X^n)\right] \leq (|\mathcal{X}|-1)\log(n+1)$

$$+ n \mathbb{E}_{p^n}\left[H\left(P_{X^n}\right)\right]$$

(by Jensen's ineq. since $H(P)$ is concave in $P$)

$$\leq (|\mathcal{X}|-1)\log(n+1)$$
$$+ n H\left(\underbrace{\mathbb{E}_{p^n}\left[P_{X^n}\right]}\right) \to P$$

$\mathbb{E}_{p^n}\{P_{X^n}(a)\} = \dfrac{n}{n}$

$$= (|\mathcal{X}|-1)\log(n+1) + n H(P)$$

Thus, $\quad \mathbb{E}_{p^n}\left[l(X^n)\right] - n H(P) \leq (|\mathcal{X}|-1)\log(n+1)$

$$\Rightarrow \boxed{ r_n^* \leq (|\mathcal{X}|-1)\log n }$$

[C] Arithmetic Code: An online scheme

$Q^{(n)}(x^n)$ pmf over sequences of length $n$ from $\mathcal{X}$

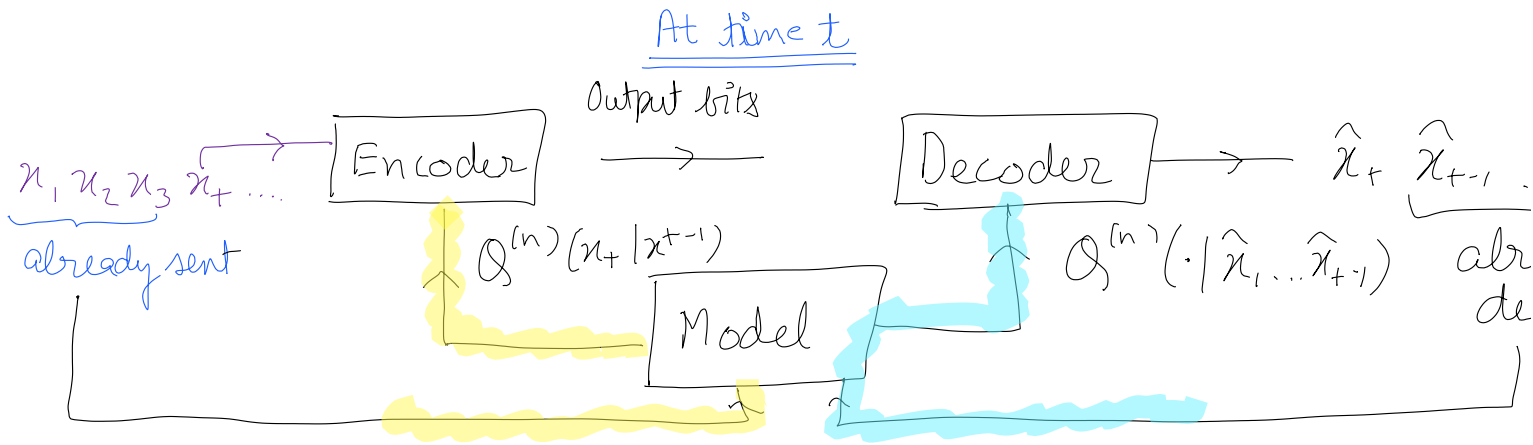$\to$ We want a prefix-free code with codeword lengths $\lceil \log$

$\hookrightarrow$ a scheme that uses $Q^{(n)}(x_i | x^{i-1})$ to encode symbol after it has encoded $(x_1, \ldots, x_{i-1})$ in the previous slots

$\hookrightarrow$ "Streaming implementation", can operate symbol-by-sy

$\hookrightarrow$ at the decoder side, upon decoding $\hat{x}$

At time t

Output bits →

$x_1, x_2, x_3, x_t, \ldots$ → [Encoder] → [Decoder] → $\hat{x}_t \; \hat{x}_{t-1}$

already sent

$Q^{(n)}(x_t | x^{t-1})$

[Model]

$Q^{(n)}(\cdot | \hat{x}_1 \ldots \hat{x}_{t-1})$ alr de

## Arithmetic code (interval representation)  Maintain an inte... $I_{t-1}$

$(x_t$

$[\ \ \ \ \ \ \ \ \ )$

(1) **Initial step**: $I_0 = [0, 1)$

(2) At time $t$,

(a) Divide $I_{t-1}$ into sub-intervals $\{I_{t-1, x}, x \in \mathcal{X}\}$

$(|I| \equiv$ length of interval $I)$

s.t. $\dfrac{|I_{t-1, x}|}{|I_{t-1}|} = Q^{(n)}\left(X_t = x \mid X^{t-1} = x^{t-1}\right)$

(b) Let $I_t = I_{t-1, x_t}$

(3) When you stop, you have an interval $I_n$ of length

$$|I_n| = |I_1| \cdot \frac{|I_2|}{|I_1|} \cdots \frac{|I_n|}{|I_{n-1}|} = Q^{(n)}(x_1) \, Q^{(n)}(x_2 | x_1) \cdots$$

$$= Q^{(n)}(x_1, \ldots, x_n)$$

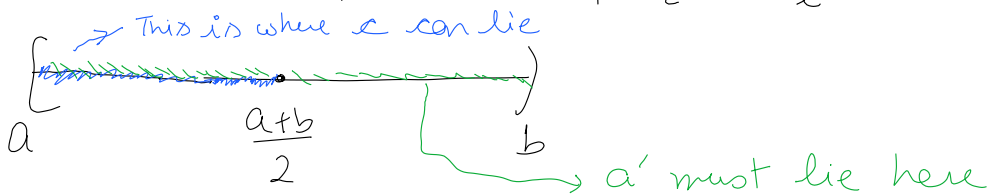How do we represent these final intervals using a prefix-free code?

## Shannon-Fano-Elias code

Represent the interval $I = [a, b)$ using the $\boxed{\ell = \lceil -\log(b-a) \rceil + 1}$
most significant bits in the binary representation of $\boxed{\dfrac{a+b}{2}}$

$\boxed{\text{This code is prefix free}}$   Suppose $\dfrac{a+b}{2} = 0.z_1 z_2 \ldots$

Let $c = 0.z_1 z_2 \ldots z_\ell$

→ This is where $c$ can lie

$[\ \ \ \ \ \ \ \ \ \ \ )$

$a \qquad \dfrac{a+b}{2} \qquad b$   → $a'$ must lie here

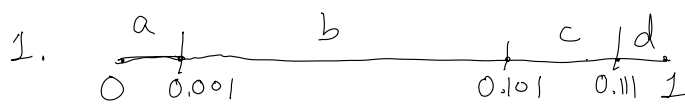Since $\ell = \lceil -\log(b-a)\rceil + 1$, $2^{-\ell} \le \dfrac{b-a}{2}$.

Suppose $y_1 \dots y_m$ has $z_1 \dots z_\ell$ as its prefix. Consider
$a' = 0.y_1 \dots y_m$. $\quad (c \le a' \le c + 2^{-\ell})$

$\Rightarrow$ Both $c$ and $a'$ must lie in $I$.

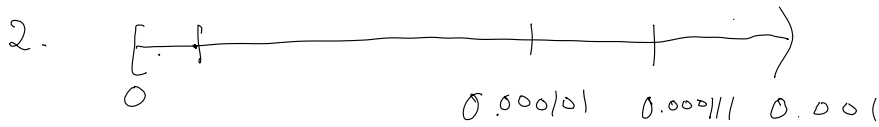$\Rightarrow$ The code is prefix-free.

$\rightarrow$ The average length of the resulting code is $\le H(P) + 1 + 1 = H($

Example: $\mathcal{X} = \{a, b, c, d\}$, $n = 4$, $Q^{(n)}(x_1, x_2, x_3, x_4) = \prod_{i=1}^{4} P(x_i)$
where $P(a) = \dfrac{1}{8}$, $P(b) = \dfrac{1}{2}$, $P(c)$

1. 

$I_1 = [0, 0.001)$

Encode the sequence $a\,c\,b\dots$

2. 

$0.000101 \quad 0.000111 \quad 0.001$

$I_2 = [0.000101, 0.000111)$

3. $\quad I_3 = [0.0001010, 0.0001010111)$

4. $\quad I_4 = [0.0001010100, 0.0001010101111)$
$\quad \underset{a}{\phantom{.}} \qquad \underset{b}{\phantom{.}}$

$\quad\hookrightarrow \dfrac{a+b}{2} = 0.00010101110$

$\quad b - a = 2^{-9} \quad \Rightarrow \quad \ell = 11$

$\Rightarrow$ The codeword corresponding to this interval is $\boxed{0001010\dots}$

## Decoder of arithmetic code
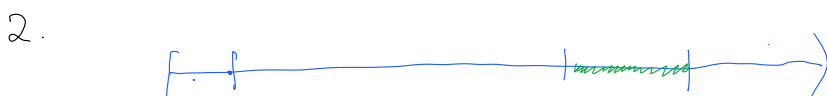
$\rightarrow$ Decoder just "inverts" the encoding procedure.
We illustrate using the example above.

Given $000101011 00$

1. 

$\Rightarrow x_1 = a$

2. 

0.00101    0.000111   0.001

$\Rightarrow x_2 = c.$
$x_3 = b, x_4 = b.$

In conclusion, we have a "practical" algorithm that attains average length $\leq H(Q^{(n)}) + 2$ and requires the model to provide $Q^{(n)}(\cdot \mid x^{t-1})$ at time $t$.

$\rightarrow$ If $Q^{(n)}$ attains $R_n^*$, we get a scheme with regret less than

## Online probability assignment and a universal scheme

<u>input:</u> $Q^{(n)}$ on $\mathcal{X}^n$ $\left( Q^{(n)}(x_+ \mid x^{t-1}) \right)$

<u>output:</u> a code with $l(\underline{x}) = \lceil \log \frac{1}{Q^{(n)}(x)} \rceil + 1$

Therefore, given a probability assignment $Q^{(n)}(\underline{x})$, arithmetic code gives a scheme with average length

$$\leq \sum_{\underline{x}} P^n(\underline{x}) \log \frac{1}{Q^{(n)}(\underline{x})} + 2$$

$$= \boxed{D(P^n \| Q^{(n)})} + n H(P) + 2$$

$\hookrightarrow$ we want to choose a $Q^{(n)}$ that minimiz
$$\max_{P \in \mathcal{P}(\mathcal{X})} D(P^n \| Q^{(n)})$$

$\rightarrow$ Further, we should be able to efficiently compute $Q^{(n)}(x_+ \mid \underline{x^{t-1}})$
(ideally, we should not be required $\hookleftarrow$
to remember the entire sequence $x^{t-1}$)

$\boxed{\text{Restrict to } \mathcal{X} = \{0, 1\}}$ We use a "Bayesian heuristic" and assum the sequence $x^n$ was generated from iid $\mathrm{Ber}(p)$, where $p \sim \mathrm{unif}$ (
(i.e., $Q^{(n)}(\underline{x})$ is given by first generating $p \sim \mathrm{unif}[0,1]$ and then taking $n$ indep. samples from $\mathrm{Ber}(p)$).

$$Q\,(X_{t+1}=1\mid X_1=x_1,\,X_2=x_2,\dots,X_t=x_t) = Q\cdot(X_{t+1}=1,\;X_j=x_j,\;\dots$$

$$\overline{Q^{(n)}(X_j=x_j,\;1\le\dots}$$

$$= \frac{\displaystyle\int_0^1 p^{k+1}(1-p)^{t-k}\,dp}{\displaystyle\int_0^1 p^{k}(1-p)^{t-k}\,dp}\,,\quad\text{where }k=\sum_{j=1}^t x_j.$$

Note that
$$\int_0^1 p^m(1-p)^n\,dp = \frac{n}{m+1}\int_0^1 p^{m+1}(1-p)^{n-1}\,dp$$

$$\vdots$$

$$= \frac{n}{m+1}\cdot\frac{(n-1)}{m+2}\cdot\;\dots\;\frac{1}{m+n}\int_0^1 p^{m+n}\,dp$$

$$= \frac{n!\,m!}{(n+m)!}\cdot\frac{1}{(n+m+1)}$$

$\longrightarrow$ empirical esti
aftu adding
each

Thus, $\quad Q^{(n)}(X_{t+1}=1\mid X^t=x^t) = \boxed{\dfrac{k+1}{t+2}} = \dfrac{k+1}{(k+1)+(t-k+1)}$

$\underset{\text{\# of 1s in }x^t}{\underbrace{\qquad}}$ $\qquad$ $\underset{\hookrightarrow\;\text{\# of 0s in }x^t}{\underbrace{\qquad\qquad}}$

This probability assignment is called the $\boxed{\text{add}-1\text{ estimate}}$ $\longrightarrow$ L

In general, we can consider the $\underline{\text{add}-\alpha}$ estimate, which corr
to a different prior on $p$.

The add$-\tfrac{1}{2}$ is known to ha
the "best" performance, and it
to Jeffrey's prior $\left(\pi(p)\propto\;\dots\right.$

$\longrightarrow$ Consider a sequence $x^n$ with $k$ ones (and $(n-k)$ zeros).

e.g.
$$Q\,(0010) = Q(0)\cdot Q(0\mid0)\,Q(1\mid00)\,Q(0\mid001)$$
$$= \frac{1}{2}\cdot\frac{2}{3}\cdot\frac{1}{4}\cdot\frac{3}{5}$$

$$Q\,(0100) = \frac{1}{2}\cdot\frac{1}{3}\cdot\frac{2}{4}\cdot\frac{3}{5}$$

In general, $\quad Q^{(n)}(x^n) = \dfrac{(1\cdot2\dots k)(1\cdot2\cdot\dots\cdot(n-k))}{2\cdot3\cdot\dots\cdot(n+1)} = \dfrac{1}{(n+1)}\cdot\dfrac{1}{\binom{n}{k}}$, where $k$ i

\# of ones in $x^n$.

$$D(P^n \| Q^{(n)}) = \sum_{\underline{x}} P^n(\underline{x}) \log \frac{1}{Q^{(n)}(\underline{x})} - n H(P)$$

$$= \mathbb{E}_{P^n} \left[ \log (n+1) \binom{n}{k} \right] - n H(P)$$

$$\leq \log(n+1) + n \mathbb{E}_{P^n} \left[ h \left( \frac{k}{n} \right) \right] - n H(P)$$

$$\leq \log(n+1).$$