

"Numerical" random variables

If $X \sim P$ is RV with values in $\mathcal{X} \subseteq \mathbb{R}$:

Expectation value (mean): $EX = E[X] = \sum_x P(x) \cdot x$

* $E[f(x)] = \sum_x P(x) \cdot f(x)$ "law of the unconscious statistician"

* $E[cX] = c \cdot E[X]$ & $E[X+Y] = E[X] + E[Y]$ (A)

* If X, Y independent: $E[XY] = E[X] \cdot E[Y]$ $\sum_{x,y} P(x)P(y)xy$

↳ $X \sim \text{Uniform}(\{-1, 1\})$, $Y = -X \Rightarrow E[XY] = -1$, $E[X] = E[Y] = 0$
NOT indep

Variance: $\text{Var}(X) = E[(X - EX)^2]$
 $= \sum_x P(x)(x - EX)^2 = E[X^2] - E[X]^2$

* $\text{Var}(cX) = c^2 \text{Var}(X)$

* If X, Y independent:

$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$ (B)

↳ we that $E[XY] = E[X] \cdot E[Y]$

Examples

P	Bernoulli(f)	Binomial(n, f)
E	f	$n \cdot f$ (A)
Var	$f(1-f)$	$n \cdot f \cdot (1-f)$ (B)

$$E[(X - EX)^2] = E[(X - f)^2] = f(1-f)^2 + (1-f)(0-f)^2 = f(1-f)$$

Three results that give these meaning:

Markov inequality: If $X \geq 0$: $\Pr(X \geq t) \leq \frac{E[X]}{t} \quad (\forall t > 0)$

Pf: $\Pr(X \geq t) = \sum_{x \geq t} P(x) \leq \sum_{x \geq t} P(x) \frac{x}{t} \leq \frac{E[X]}{t} \quad \square$

Chebyshev inequality: $\Pr(|X - EX| \geq \varepsilon) \leq \frac{\text{Var}(X)}{\varepsilon^2}$

With high probability (w.h.p) deviation from mean is of order $\sqrt{\text{Var}(X)}$

Pf: Apply Markov to $Y = (X - EX)^2$. \square

Law of large numbers: Suppose $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P$ with $\begin{cases} \text{mean } \mu, \\ \text{variance } \sigma^2. \end{cases}$

Let $\bar{X} := \frac{1}{n} (X_1 + \dots + X_n)$. Then:

$$\Pr(|\bar{X} - \mu| \geq \varepsilon) \leq \frac{1}{n} \frac{\sigma^2}{\varepsilon^2}$$

WHP: empirical average \approx expectation value

PF: $E\bar{X} = \mu$ & $\text{Var}(\bar{X}) = \frac{1}{n^2} \text{Var}(X_1 + \dots + X_n) = \frac{\sigma^2}{n}$. \leadsto Chebyshev. \square

Convex and concave functions (§2.7)

Suppose $f: I \rightarrow \mathbb{R}$ is function on interval $I = (a, b)$ $a = -\infty$ or $b = \infty$ allowed

We say f is **convex** if $f'' \geq 0$

 \exp, x^2, \dots

Concave if $f'' \leq 0$

 \log, \sqrt{x}, \dots

Jensen's inequality: Let Z be a RV.

$$\begin{aligned} \text{If } f \text{ is convex: } & E[f(Z)] \geq f(EZ) \\ \text{If } f \text{ is concave: } & E[f(Z)] \leq f(EZ) \end{aligned}$$

$$\text{i.e. } \sum_z P(z) f(z) \begin{cases} \geq \\ \leq \end{cases} f\left(\sum_z P(z) z\right)$$

If $f'' > 0$ or $f'' < 0$: "=" holds only if Z is constant ∇

Entropy (§2.4)

Entropy of a random variable (RV) X with distribution P :

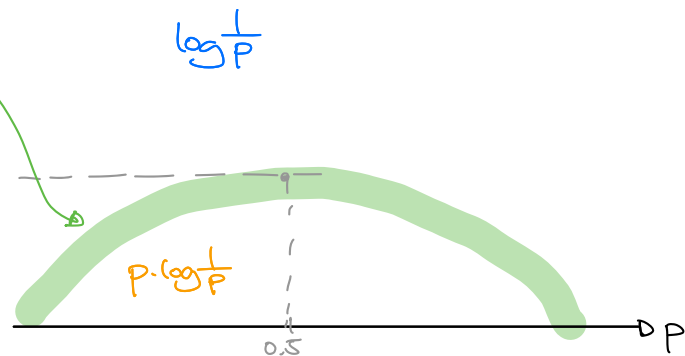
$$H(X) := H(P) := \sum_x P(x) \cdot \log \frac{1}{P(x)} = E\left[\log \frac{1}{P(X)}\right] \quad \text{unit "bit"}$$

$0 \cdot \log \frac{1}{0} = 0$ always base 2

eg. $X \sim \text{Bernoulli}(p)$: **binary entropy**

$$H(X) = p \cdot \log \frac{1}{p} + (1-p) \cdot \log \frac{1}{1-p}$$

$e \in [0, 1]$



Properties:

* $H(X) \geq 0$, = iff constant $p \cdot \log \frac{1}{p} \geq 0 \quad \forall p \in [0,1], = \text{iff } p=0 \text{ or } p=1$

* $H(X) \leq \log \#\{x: P(x) > 0\} \leq \log \#\Omega_X$
 $H(X) = \log \#\Omega_X \iff X \text{ uniformly random}$

Pf: Apply Jensen with $f = \log$ and $Z = \frac{1}{P(X)}$:
 $E[\log \frac{1}{P(X)}] \leq \log E[\frac{1}{P(X)}]$
 with equality iff $P(X)$ constant, i.e.
 $P(x) > 0, P(y) > 0 \implies P(x) = P(y) \quad \square$

* NOTATION: $H(X, Y) = H(XY)$ = entropy of joint distribution $P(x, y)$

If X, Y independent: $H(X, Y) = H(X) + H(Y)$

Pf: Since $P(x, y) = P(x)P(y)$ we have $\log \frac{1}{P(x, y)} = \log \frac{1}{P(x)} + \log \frac{1}{P(y)}$
 \implies take expectation values. \square

Interpretation? Let us call $h(x) = h(X=x) = \log_2 \frac{1}{P(x)}$ the information content (or "surprisal") of an outcome $x \in \Omega_X$.

$\implies H(X) = E[h(X)]$ is average information content.

Why is this a good definition? Three suggestive examples:

① Uniformly random number in $\{0, \dots, 255\}$: $H(X) = \log_2 256 = 8 \text{ bit}$

A							
B							
C							
D							
E							
F							
G			⊗				
H							
	1	2	3	4	5	6	7

② Poor man's submarine game: Single submarine hidden, other player asks if submarine in some square \rightarrow hit/miss

1st move: $P(\text{hit}) = \frac{1}{64} \rightarrow h(\text{hit}) = 6 \text{ bit}$ learned precise location (64 options)
 $P(\text{miss}) = \frac{63}{64} \rightarrow h(\text{miss}) \approx 0.0221 \text{ bit}$ learned little (63 remaining)

③ "Wenglish" has 2^{15} words in $\{A, \dots, Z\}^5$ s.t. frequency of single letters matches English. Let w be uniformly random word in this list.

$H(w) = 15 \text{ bit}$, i.e. on average 3 bit/letter

but e.g. $p(w_1 = z) = 0.1\% \implies h(w_1 = z) \approx 10 \text{ bit}$ \leftarrow no contradiction: we learn less info from the rest since few words start with z