

Last time:

* Approximations + self consistency

$$\epsilon x^5 + x - 1 = 0 \Rightarrow x \approx 1 \quad (\epsilon x^2 \approx \epsilon \ll 1)$$

↑ ↑ ↑
small? dominant?

(vs $x \approx \epsilon^{-1/5}$ when $\epsilon \gg 1$)

* Probability ($x \sim p(x)$)

\Rightarrow Generating functions: $H_x(z) \equiv \langle e^{-zx} \rangle = \int e^{-zx} p(x) dx$

e.g. Poisson $p(n) = \frac{\lambda^n}{n!} e^{-\lambda} \Leftrightarrow H_n(z) = e^{-\lambda(1-e^{-z})}$

\Rightarrow Central limit theorem: $X_1, X_2, \dots, X_n \sim p(x)$

as $n \rightarrow \infty \Rightarrow \frac{1}{n} \sum_{i=1}^n X_i \rightarrow \text{Gaussian}(\langle x \rangle, \frac{\text{Var}(x)}{n})$

Today:

- ① Intuition about probability
- ② Biological background (#s + scales)
- ③ Simple model of evolution (if time permits)

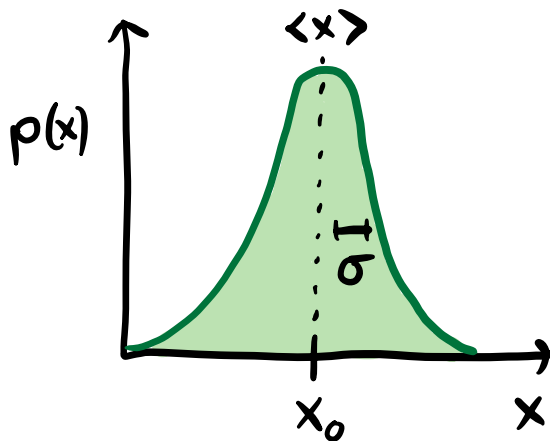
Intuition about random variables ("average" vs "typical")

Probability can be hard b.c. it forces us to think about many outcomes all @ once...

⇒ often want some way of summarizing "typical" behavior

we will frequently encounter 2 broad classes:

Case 1 ("fuzzy noise"):

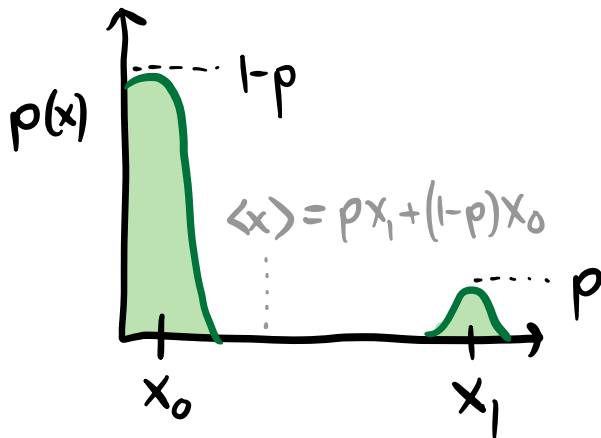


e.g. Binomial(N, p)
when $Np \gg 1$ + $N(1-p) \ll 1$

$$\Rightarrow X \approx X_0 \pm \sigma$$

⇒ average is good summary of "typical"

Case 2 ("jagged noise"):



e.g. Binomial(N, p)
when $Np \ll 1$
e.g. did a mutation occur or not?

\Rightarrow no realization of x has $x \approx \langle x \rangle$

\Rightarrow average is poor summary of typical!

(better off guessing $x \approx x_0$, + rare exceptions)

\Rightarrow distinction becomes important if we do something w/ x :

e.g. $y = F(x) =$ "future growth of x mutations"

\Rightarrow in case 1: can treat noise as small perturbation
& use approx. methods above.

e.g. using Taylor expansion around $x \approx x_0$:

$$F(x) = F(x_0 + (x - x_0)) \approx F(x_0) + F'(x_0)(x - x_0)$$

w/ $x = x_0 \pm \sigma$ \downarrow

$$y \approx F(x_0) \pm F'(x_0)\sigma$$

\uparrow
deterministic
guess

\nwarrow small spread
("fuzziness")
due to noise.

case 2:

$$y = \begin{cases} F(x_0) \\ F(x_1) \end{cases}$$

w/ prob $1-p$

\swarrow this can be
"typical" case
(most of time)

w/ prob p

\nwarrow "rare event"
happens repeatedly.

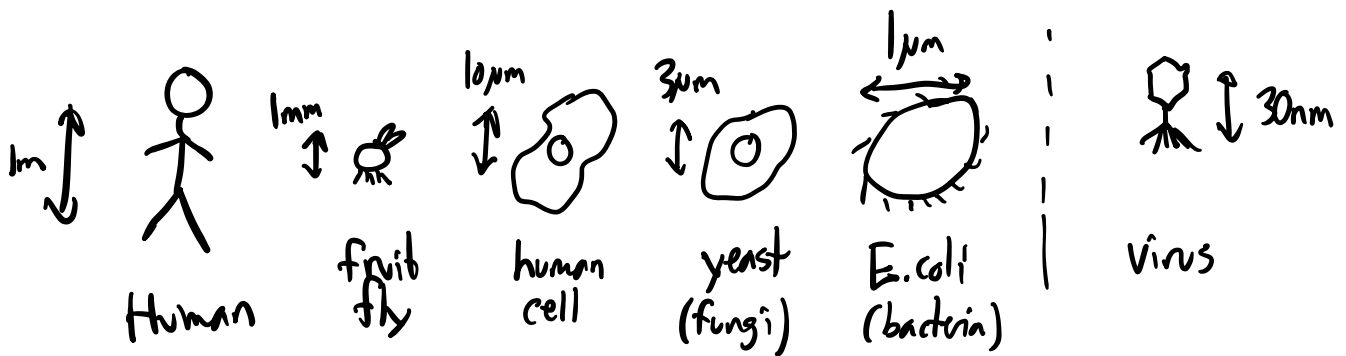
\Rightarrow often used to thinking about case 1 dist'n s.

\Rightarrow evolution will have many examples of case 2!

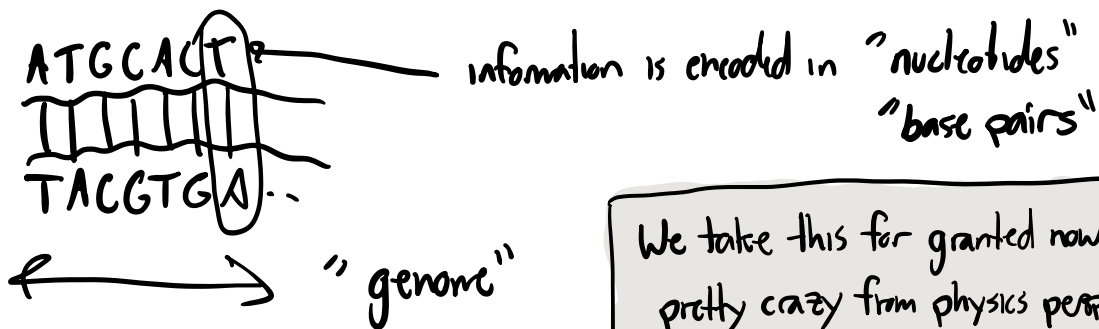
Biological background (key #'s and scales)

① Organisms come in huge range of shapes + sizes:

"Model organisms" we will encounter in this course:



② Despite diffs, these organisms are similar in that instructions to create them are encoded in a single* long molecule of DNA:



We take this for granted now, but pretty crazy from physics perspective! (important info in 1 molecule vs many)

$4^3 = 64$ different codons \rightarrow 20 amino acids
+ "start codon"
+ "stop" codon
"genetic code"
 \Rightarrow has degeneracy

\Rightarrow typical protein ~ 300 AA (1000bp of DNA)

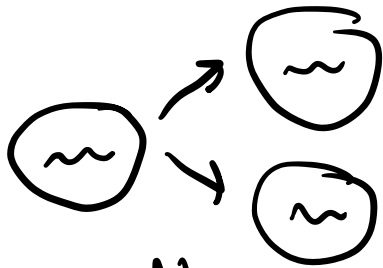
\Rightarrow # of genes varies widely across organisms:

humans: 20,000 genes yeast: 6,000 genes
E. coli \sim 4,000 genes viruses \sim 10 genes.

\rightarrow 1000x bigger genome \Rightarrow but 5x as many genes.

\Rightarrow rest of genome is "noncoding" \rightarrow regulation
("coding" = genes) \rightarrow "junk"

\Rightarrow net effect of doing all these things
is that the organism makes a copy of itself:



Δt
 ← →
 | doubling time
 | "generation"

- ① new cell wall, all other proteins (including ribosomes!)
- ② needs to copy its DNA (DNA polymerase) (not usually limiting factor in growth)

Some characteristic generation times:

humans: ~ 20 yrs

E. coli ~ 20 mins - 1 hr (lab)

human cells: ~ 1 day (HeLa)

1 hr - 1 day? (in gut)

Prochlorococcus ~ 1 day

Virus: HIV ~ 15 hrs
SARS-CoV-2 ~ 10 hrs

(ocean bacterium, one of the most abundant photosynthetic organisms on earth, $N \sim 10^{27}$)

⇒ Since $n=1$ genome, can make errors during copying

... ATGCCA ... parent
... ATG↓TCA ... offspring

"mutations"

⇒ simplest mutations are "point mutations" ($A \rightarrow T, T \rightarrow C, \dots$)
aka "single nucleotide mutations" / "substitutions" / "SNPs"

⇒ can also have "insertions":

... ATGTTTCA ...
↓
... ATGTTTTC ...
(+3T)

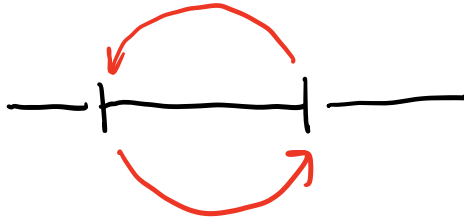
⇒ or "deletions"

... ATGTTTCA ...
↓
... ATGTCA ...
(-2)

e.g. slippage of DNA pol

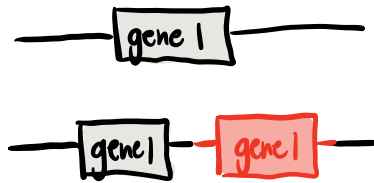
⇒ can also have larger "structural rearrangements":
(eg. > 1kb)

e.g. "inversions"

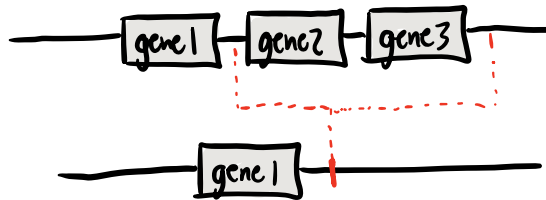


often mediated
by special genes
known as
"transposons"

e.g. "duplications"



e.g. "deletions"



⇒ upshot: can get pretty complicated

⇒ cells have sophisticated machinery
for fixing errors that occur in genome...

⇒ net mutation rates (μ) vary across organisms!

e.g. Humans: $\mu \sim 10^{-8}$ single nucleotide muts/bp/gen

Human cells: $\mu \sim 10^{-10}$ /bp/division E. coli: $\mu \sim 10^{-10}$ /bp/gen

viruses: up to $\mu \sim 10^{-5}$ /bp/gen (SARS-CoV2 10^{-6} /bp/gen)

⇒ Using these #'s, can already make some interesting predictions...

Evolutionary "Fermi Problems"

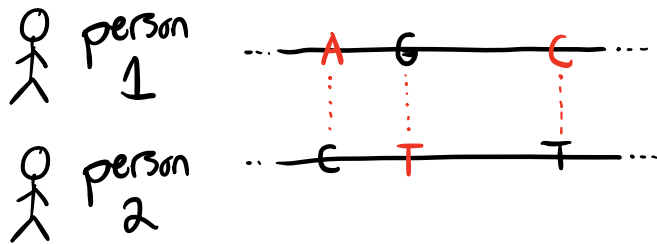
e.g. in Humans genome is $L = 3 \times 10^9$ bp (actually $\times 2$, since there are 2 copies of each chromosome "diploid")
+ mutation rate $\mu \sim 10^{-8}$ /bp/gen

$$\Rightarrow L \cdot \mu = 3 \times 10^9 \frac{\text{bp}}{\text{genome}} \times 10^{-8} \frac{\text{mutations}}{\text{bp} \cdot \text{gen}} \approx 30 \text{ mutations per genome per gen.}$$

⇒ there are $N \sim 10^{10}$ humans on earth, so

$$\Rightarrow N \times \mu \sim 10^{10} \times 10^{-8} \sim 100 \text{ mutations produced @ every site in human genome per generation (in some individual)}$$

⇒ but, if we pick 2 random people + compare genomes:



Empirically: differ @ ~0.1% of genome

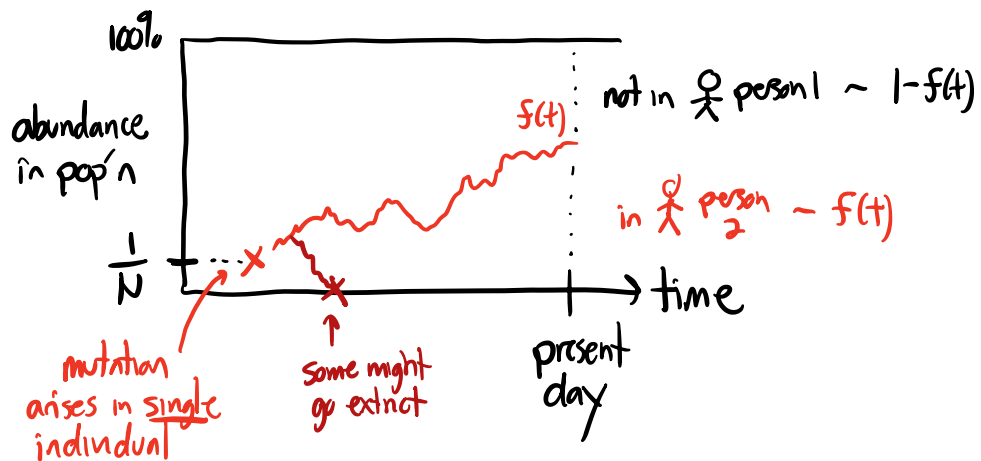
Question: what sets this scale? why not 10^4 or 10^{-2} ?

⇒ one factor: dynamics over time

⇒ for mutations produced in previous generation,

$$\Pr[\text{mut in person 1 or person 2}] \approx 100 \text{ muts/bp} \times \frac{2}{10^{10}} \sim 10^{-8} / \text{bp}$$

⇒ many differences observed today occurred in past



⇒ must understand mut'n trajectories over time ("dynamics")
(goal of next several lectures...)

Another Fermi calculation:

* all single mutations produced every gen in humans
but all pairs of mutations are not:

$$\Pr \left(\begin{array}{l} \text{site 1 + site 2} \\ \text{mutated in same} \\ \text{newborn} \end{array} \right) \sim N \times \mu \times \mu \sim 10^{10} \times 10^{-8} \times 10^{-8} \sim 10^{-6}$$

⇒ must wait $\sim 10^6$ gens (20 million yrs!) for
a given pair of sites to mutate @ same time

⇒ Upshot: past dynamics even more important
for combinations of mut'n's.

⇒ can also repeat same calculations for E. coli...

⇒ genome is $L = 4 \times 10^6$ bp + $\mu \sim 10^{-10}$ /bp/gen

$$\Rightarrow L \times \mu \sim 4 \times 10^{-4} \text{ mutations/genome/gen}$$

$\Rightarrow > 1000$ replications before a single error!

$$\Rightarrow N_g \sim 10^9 - 10^{10} \text{ E. coli cells in single person's gut}$$

$$\Rightarrow N_g \times \mu \sim 0.1 - 1 \text{ (almost every bp mutated w/in us each day)}$$

$$\Rightarrow N_h \sim 10^{10} \text{ guts in human pop'n,}$$

$$\Rightarrow N_h \times N_g \times \mu \times \mu \sim 0.1 - 1 \Rightarrow \text{almost all double mutations produced in worldwide E. coli pop'n each day}$$

$$\Rightarrow \text{but not triple mutants } (10^{10} \times 10^{10} \times (10^{-10})^3 \ll 1)$$

\Rightarrow more generally, for single gene of $L \sim 1000$ bp

$$\Rightarrow 4^L \approx 10^{600} \text{ possible DNA sequences!}$$

compare to $\sim 10^{82}$ atoms in universe

\Rightarrow sequence space is very big (& sparsely populated)

What do mutations do? "genotype \Rightarrow phenotype map"

\Rightarrow in general, we don't know a priori (even for model organisms like E. coli!)

\Rightarrow but in special cases, can make some guesses based on structure of the genetic code...

e.g. if mutation occurs in a gene:

\Rightarrow changes a codon (e.g. ATC \rightarrow ATT)

① due to degeneracy, codon could code for same AA

\Rightarrow doesn't change protein "synonymous mutation"

② could change to something else "nonsynonymous mut'n"

\hookrightarrow e.g. other AA (small change?) "missense mut'n"

\hookrightarrow e.g. stop codon \Rightarrow truncates gene (big change)

"loss-of-function" / "nonsense" mut'n