So far, we've talked about systems that "just"
react. Even with systems that have memory,
we can't say they're learning. Because they
never change their policies. Even Bayesian
networks don't really "learn" as much as
process sensory data with uff. conditions.

That's because we need some way of
assigning value to an action. Was this
the right or wrong choice? Does this work
or not work? This is called UTILITY.

utility, as a concept, comes from economics.
Often, it's talked about abstractly, ie. it
has no unit. But concretely it's often
operationalized as price.



you         I        → we → value =
have       want        negotiate.   $ price
chocolate    it                     =
bar                                  $1.00

so we can say the chocolate bar
is worth $1.00 and must therefore
give me $1.00 of utility.

I want to print out a divergence between this model and reality.

1. Humans can predict, recall, learn. $1 ≠ $1.00 if your $1/$1000 vs. $1/1000000.

2. Most things don't have $ value. "Money can't buy you love" is literally true.

3. Even within shift preference models, price shifts. E.g. Kahneman.

But utility is (potentially) computable. So, it is the basis of our modelling.

Blacksmith in Dist. Kingdom.

Purse $\frac{lnw:}{raw}$ Labour theory of value

10 coins

prod = +1 to value and Inventory.

1 unit of raw = $1.

King ^Bezos ^von. has set up a distribution network. Each day, a cart will ship. You can:

1) ~~gross~~ sell what you've make.

2) Buy a ~~product~~ ^raw. material.

3) wait and do nothing

4) spend the day producing something.

you tell the cart man every day for tomorrow.

| Buy | sell | wait | produce. |
|-----|------|------|----------|
| -1  | +2   | $\emptyset$ | $\emptyset$ |

obviously:

sell > wait = produce > buy    = profs

If your policy is strictly argmax (profs),
you would never make or buy anything!

Instead:

| coin? | prod? | raw? | action. |
|-------|-------|------|---------|
| Y     | N     | N    | buy     |

☆ fill out the rest.

Let's say the cnt is now true random.
You probably need a more complex
policy. eg.
☆ coin < 2      coin > 10 ... etc.

and you need to kick meter it "washed"!

Now let's say you need to schedule ppl.
produce now is -10 and needs advance
commitment. But Sell is now
+40 because of efficiency.
cnt still random.
☆ schedule policy.

Now let's say that the cnts "look" random, but there actually is a pattern. The pattern is mostly consistent but not perfectly. How do you learn the best policy?

Introduce Q-Learning, a type of reinforcement learning.

$$Q : S \times A \to \mathbb{R} \quad \longrightarrow \text{reward}$$

$$\downarrow$$

$$\longrightarrow \text{action } buy, sell, wait, produce.$$

state

<u>coin</u>    <u>inv. raw</u>    <u>inv. prod.</u>

The ideas is  ~~how~~  to make a big table of time reward so:

$$\frac{c}{10} \quad \frac{raw}{0} \quad \frac{prod}{1} \quad \times \text{ sell } \to \overset{\$}{+} 40$$

may be true in a moment, but might be long-run costly if faced with a choice to buy raw. why?

$$\frac{c}{109} \quad \frac{r}{0} \quad \frac{prod}{0} \quad \cancel{\times} \text{ sell } \to \$0$$

so abstract <u>utility</u> s.t.

$$\frac{c}{19} \quad \frac{r}{0} \quad \frac{prod}{1} \quad \times \text{ sell } < \frac{c}{10} \quad \frac{r}{0} \quad \frac{prod}{1} \quad \times \text{ buy}$$

★ Time table so that that a 5-day
week can be "planned" using
just utility values.

This is hard. Essentially impossible to design
by hand w/o many restricted circumstances.

★ Translate to the robot. How could
you reward + punish your robot?

# Distribution 02

Warm up: Draw networks with
"islands"... what
does it take to make
an island into a...
"penninsula"?



vs.

Meet the profs 2 mile @ 5pm
Koerner's.

$$\underline{\underline{utility}} = value$$

$$u(action) = value$$

 ~     $1.00

$$\text{⌂} = \$1.00 = 🪙$$

equal pref.

1. $1/1000    $1/1 000 000

2. "money can't buy you love."

3. Kanheman

Blacksmith     iron → products

purse



raw

-1        +2        0         0
buy      sell      wait     priduce

products

I cart / day

actions

sell > wait = produce > buy

✗ argmax (actions) → sell

| coin? | prod? | raw? | → action ⌄ |
|-------|-------|------|------------|
| Y | N | N | buy |
| Y | Y | N | → sell |
| Y | N | Y | → produce |
| N | Y | N | → sell |
| N | N | N | → wait |

policy ↑↓

what if cart is random?

$c < 2$    vs.    $c \geq 10$

make a new policy table.

$c < 2 \rightarrow$ sell        $c > 10 \rightarrow$ prod.
$p? = y$            $p > 20$    $-5^t$ prod

for the next week, you hire
labourers. 1 iron → 10 prod.

-5 because you're paying to produce
but you get +10 products.

|  |  |  | $\frac{c}{10}$ | $\frac{raw}{5}$ | $\frac{prod}{20}$ |
|---|---|---|---|---|---|
| sell | +10 | X |  |  |  |
| buy | ∓1 |  |  |  |  |
| produce | -5 |  |  |  |  |
| wait | 0 |  | $\underline{\underline{100}}$ | $\underline{\underline{0 \div}}$ | $\underline{\underline{0 =}}$ |

↳ what else can you track
   it cnts aren't actually
   random but just
          appear so?

| day | c | raw | prod. | utility | action |
|---|---|---|---|---|---|
| 7ᵗʰ Nov | 10 | 5 | 25 | ⇒ $10 | sell |
| ⋮ |  |  |  |  |  |

Reinforcement Learning
Q - learning

$Q : S \times A \longrightarrow \mathbb{R}$   6

state

| coin | value | prod. | | actions | → Reward |
|------|-------|-------|---|---------|----------|
| 1    | 0     | 0     | | sell    |          |
| 0    | 1     | 0     | | buy     |          |
| 6    | 0     | 1     | |         |          |

buy → -1 Reward

buy, produce, sell → (+10)

1. State tracking?
2. Actions?
3. Reward?