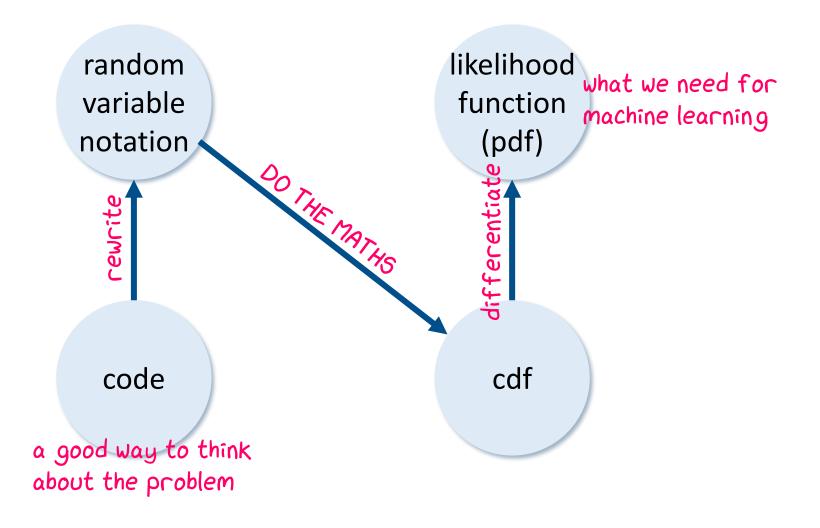
LAST LECTURE

Bespoke probability distributions: from code to likelihood

There are four ways to specify a distribution:



EXERCISE

What's the pdf for this random variable?

```
def rx(u,v,w,p):
   \# preconditions: u < v < w, and 0 
   k = np.random.choice(["left", "right"], [p, 1-p])
   if k == "left":
       return np.random.uniform(u,v)
   else:
       return np.random.uniform(v,w)
```

1 Rewrite it in random variable notation

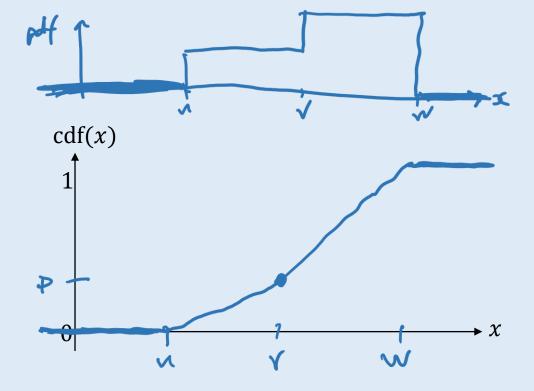
Let
$$K = \begin{cases} \text{left} & \text{with prob. } p \\ \text{right} & \text{with prob. } 1 - p \end{cases}$$
Let $X \sim \begin{cases} U[u,v] & \text{if } K = \text{left} \\ U[v,w] & \text{if } K = \text{right} \end{cases}$

2) Do the maths: simplify the cdf into elementary building blocks

$$\mathbb{P}(X \le x) = \mathbb{P}(X \le x | K = \text{left}) \times \mathbb{P}(K = \text{left}) + \mathbb{P}(X \le x | K = \text{right}) \times \mathbb{P}(K = \text{right})$$

$$= p \, \mathbb{P}(U[u, v] \le x) + (1 - p) \, \mathbb{P}(U[v, w] \le x)$$

$$= \begin{cases} \text{if } x < u: & \text{production} \\ \text{if } x \in [u, v]: & \text{production} \\ \text{if } x \in [v, w]: & \text{production} \\ \text{if } x > w: \end{cases}$$

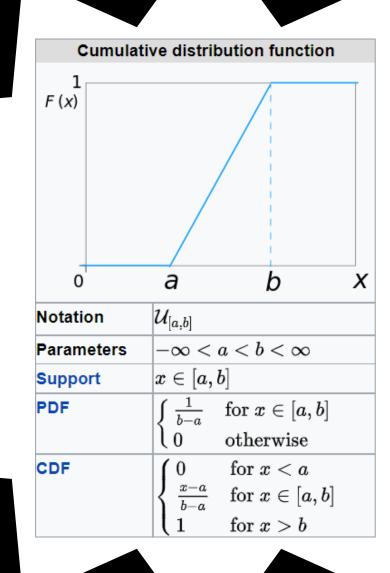


by the Law of Total Probability

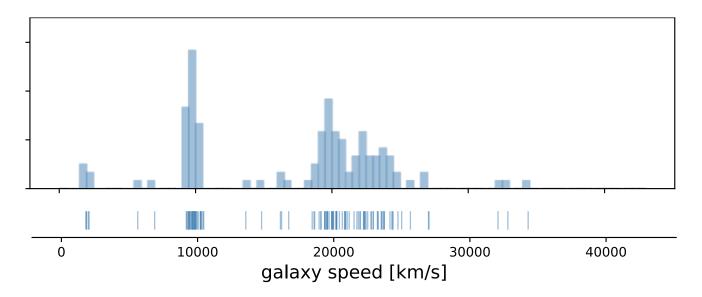
$$\mathbb{P}(X = x) = \sum_{y} \mathbb{P}(X = x | Y = y) \mathbb{P}(Y = y)$$

$$\mathbb{P}(X=x) = \sum_{y} \mathbb{P}(X=x|Y=y)\mathbb{P}(Y=y)$$
Sourily check: at $x = v$,
$$\begin{cases} P \\ P \end{cases}$$

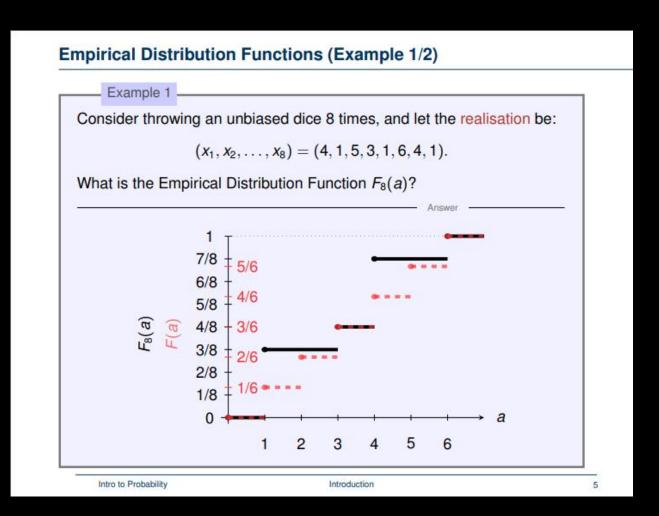
Wikipedia: Uniform distribution



Our goal: to find the best distribution we can to fit this dataset.



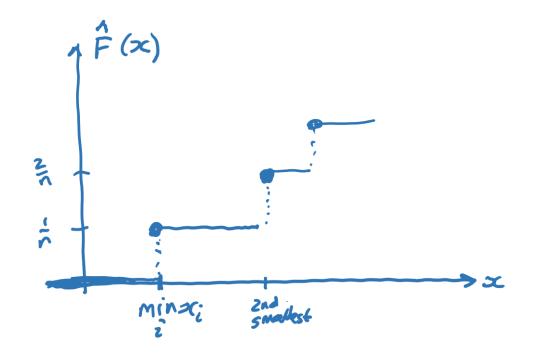
IA Probability lecture 10 Empirical cumulative distribution functions



ECDF

Given a dataset of numerical values $[x_1, x_2, ..., x_n]$, the empirical cumulative distribution function or ecdf is

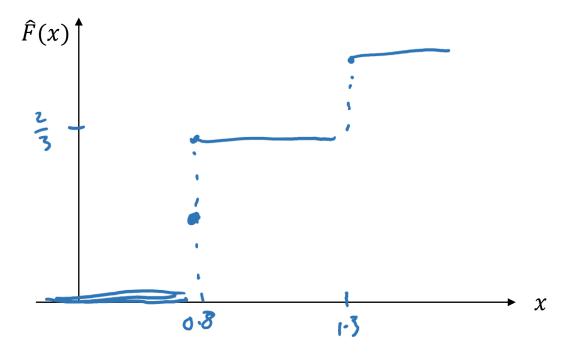
$$\widehat{F}(x) = \frac{1}{n} \begin{pmatrix} \text{how many datapoints} \\ \text{there are } \le x \end{pmatrix}$$



```
x = [...]
F = np.arange(1, len(x)+1) / len(x)
plt.plot(np.sort(x), F, drawstyle='steps-post')
```

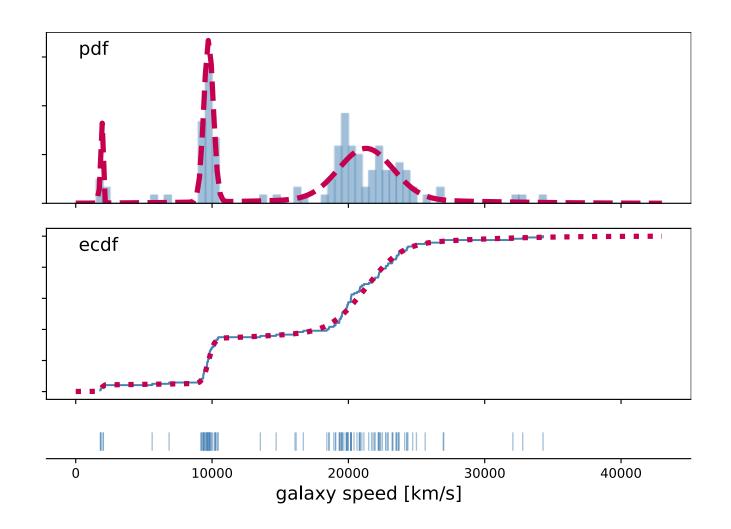
What if there are repeated values in the dataset, e.g.

$$x = [0.8, 0.8, 1.3]$$



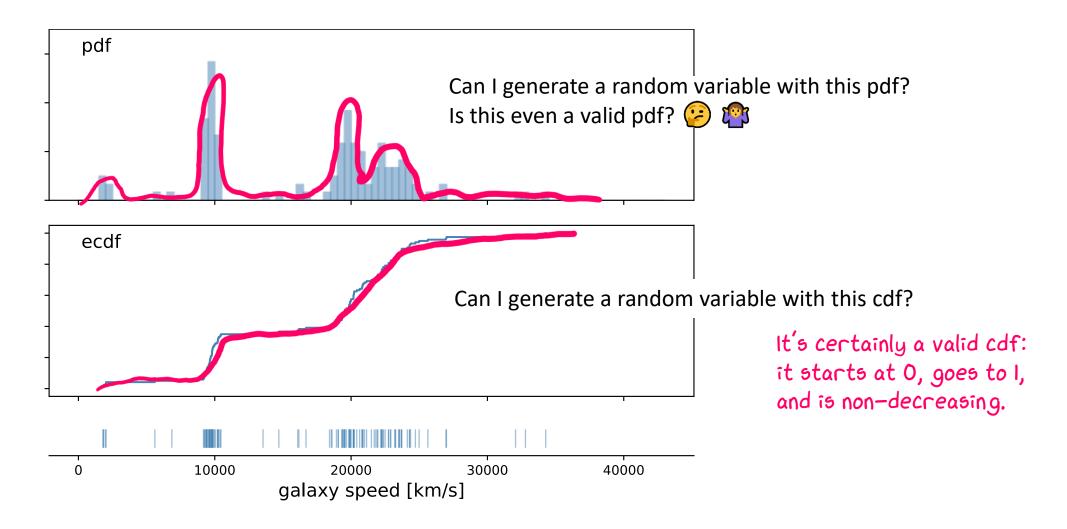
```
x = [...]
F = np.arange(1, len(x)+1) / len(x)
plt.plot(np.sort(x), F, drawstyle='steps-post')
(This code will plot an extra point at (0.8, I/3), but who cares?
```

The plot is still correct.)



fitted
---- Gaussian mixture
model

But can I find a better-fitting distribution?

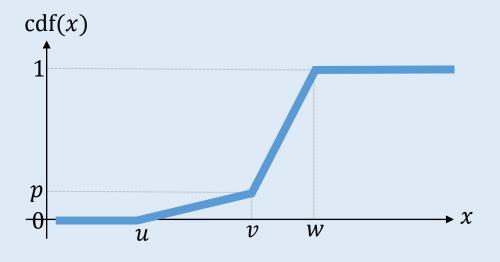


Let's build up our skills at turning cdf plots into code.

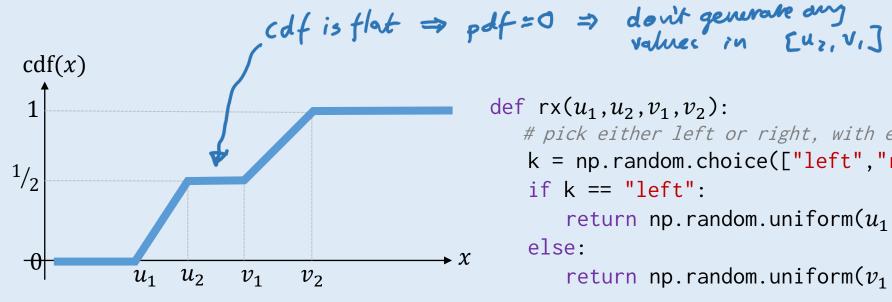
But can I find a better-fitting distribution?

This cdf has equal probabilities (p = =).

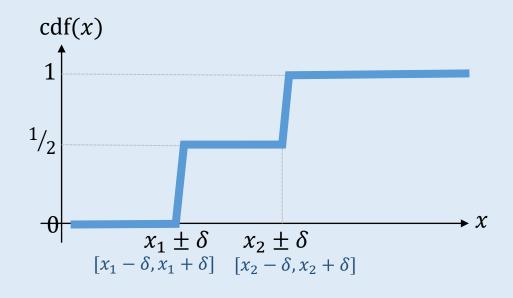
That's the defoult for /np. random.choice —



```
def rx(u, v, w, p):
   k = np.random.choice(["left", "right"], [p, 1-p])
   if k == "left":
      return np.random.uniform(u,v)
   else:
      return np.random.uniform(v,w)
```

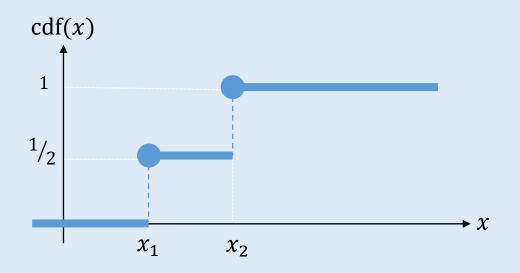


def $rx(u_1, u_2, v_1, v_2)$: # pick either left or right, with equal probability no need 6 k = np.random.choice(["left", "right"]") if k == "left": return np.random.uniform (u_1, u_2) else: return np.random.uniform (v_1, v_2)



```
\begin{aligned} &\text{def } \operatorname{rx}(x_1, x_2, \delta) \colon \\ & \text{k = np.random.choice}(["left", "right"]) \\ & \text{if } \text{k == "left"} \colon \\ & \text{return np.random.uniform}(x_1 - \delta, x_1 + \delta) \\ & \text{else:} \\ & \text{return np.random.uniform}(x_2 - \delta, x_2 + \delta) \end{aligned}
```

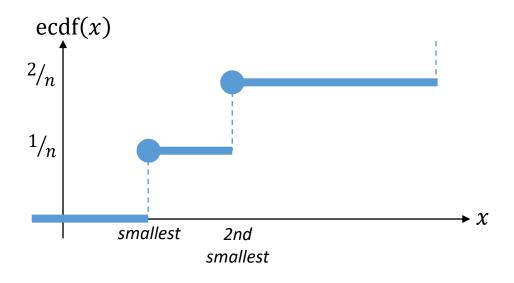
Extreme rose as F-00



def $rx(x_1, x_2)$: k = np. random. choice(["left", "kight"])if k = ||left||:

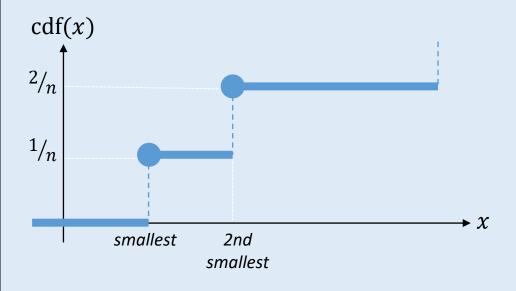
return x_1 else:

return x_2 ||left|| ||l



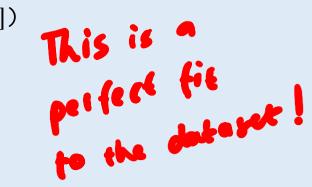
Recall the empirical distribution for a dataset $\vec{x} = (x_1, x_2, ..., x_n)$:

$$\operatorname{ecdf}(x) = \frac{1}{n}(\#\operatorname{points} \le x)$$



To generate a random variable \hat{X} whose cdf matches exactly this step function:

def rxhat(
$$[x_1,...,x_n]$$
):
return np.random.choice($[x_1,...,x_n]$)



The empirical distribution

Given a dataset $[x_1, x_2, ..., x_n]$ let \hat{X} be the random variable obtained by picking one of the x_i at random. (This is a discrete random variable.)

We say this random variable has the empirical distribution of the dataset.

The ecdf only applies to real-valued random variables, whereas this definition makes sense for any type of data (text, images, etc.)

Instead of saying "the cdf of \widehat{X} matches the ecdf of the data", we can say

$$\mathbb{P}(\hat{X} \in A) = \frac{1}{n} \sum_{i=1}^{n} 1_{x_i \in A}$$
$$\mathbb{E} h(\hat{X}) = \frac{1}{n} \sum_{i=1}^{n} h(x_i)$$

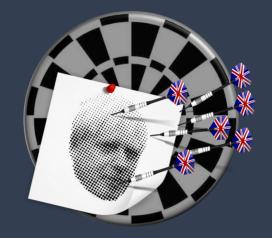
Applications of the empirical distribution

"Whenever you want to work with a true distribution, you can just bung in an empirical distribution instead."

Monte Carlo

When we want probabilities / expectations but the maths is too hard generate a sample and work with it instead.

$$\mathbb{E} h(X) \approx \frac{1}{n} \sum_{i=1}^{n} h(x_i) = \mathbb{E} h(\widehat{X})$$



Holdout approximation

When we want probabilities / expectations but we don't know the true distribution, find a sample and work with it instead.

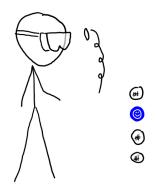
$$\mathbb{E} h(X) \approx \frac{1}{n} \sum_{i=1}^{n} h(x_i) = \mathbb{E} h(\hat{X})$$





The empirical distribution is a perfect fit for a dataset. Why bother fitting a parametric probability model at all?

§9. The frequentist approach to generalization



I tossed four coins and got x=1 head. My data model is $X \sim \text{Bin}(4, \theta)$

What can I say about the **laws of nature**, i.e. about the true value of θ ?



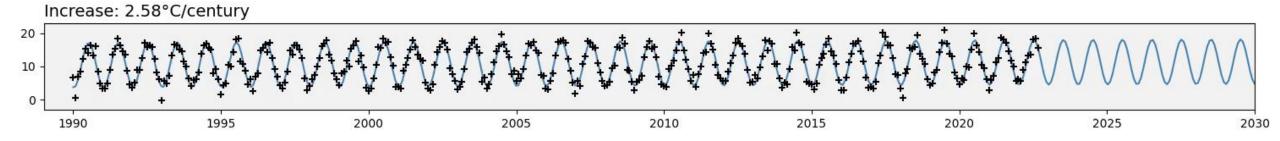
I saw X=1. Let me go figure out how likely is each possible explanation $\Theta=\theta$.

Bayes's rule:

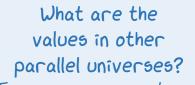
$$Pr_{\Theta}(\theta|x) = \kappa Pr_{\Theta}(\theta) Pr_X(x|\Theta = \theta)$$

I saw X=1, $\hat{\theta}$ =1/4, IN THIS REALITY. What was $\hat{\theta}$ in other dimensions of the multiverse?





I see temperatures rising by $\hat{\gamma}$ =2.58°C / century, in this reality.



<CS VERSION>

How might I simulate the parallel universes?

Climate confidence challenge.

Find a 95% confidence interval for γ , for Cambridge from 1985 to the present. (It's your choice how to simulate the multiverse.)

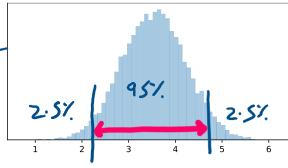
Please submit your answer on Moodle by Tuesday 11 November

Confidence intervals via resampling

Given a dataset x,

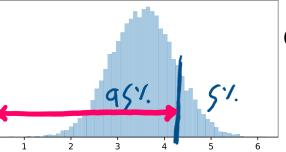
- 1. Decide on a readout function t(x)
- 2. "Simulate a multiverse of datasets."
 - Fit a model for the dataset
 - Let X* be a random synthetic dataset, generated from the fitted model
 - Simulate many synthetic datasets
- 3. Compute t for each dataset, and report the spread of t for example with a histogram or a confidence interval

This has to be a <u>computable</u> function of x, i.e. it's not allowed to have any unknown parameters. Such a function is called a statistic.



Two-sided 95% confidence interval

np.quantile(tsamples, [.025, .975])



One-sided 95% confidence interval

np.quantile(tsamples, [0,.95])

```
Example.
```

We are given a dataset

x = [4.3, 5.1, 6.1, 6.8, 7.4, 8.8, 9.9]

which we decide to model as independent samples from $N(\mu, \sigma^2)$. Find a 95% confidence interval for $\hat{\mu}$.

- 8 # 3. Sample the readout statistic, and report its spread
- 9 t_ = [t(rx_star()) for _ in range(10000)]
- 10 lo,hi = np.quantile(**t_**, [.025, .975])

```
Example 9.2.1.
```

We are given a dataset

x = [4.3, 5.1, 6.1, 6.8, 7.4, 8.8, 9.9]

lo, hi = np.quantile(t_{-} , [.025, .975])

which we decide to model as independent samples from $N(\mu, \sigma^2)$. Find a 95% confidence interval for $\hat{\mu}$.

```
# 1. Define a readout statistic
def t(x): return np.mean(x)

# 2. To generate a synthetic dataset ...

# μhat = np.mean(x)

# ohat = np.sqrt(np.mean((x-μhat)**2))

# def rx_star():
return np.random.normal(loc=μhat, scale=σhat, size=len(x))

# 3. Sample the readout statistic, and report its spread
# 4. = [t(rx_star()) for _ in range(10000)]
```

Confidence intervals via parametric resampling

Given a dataset x and a parametric probability model $Pr(x; \theta)$

- 1. Decide on a readout function t(x)
- 2. "Simulate a multiverse of datasets."
 - Fit this model, i.e. estimate $\hat{\theta}$
 - Let X* be a random synthetic dataset, generated from the fitted model
 - Simulate many synthetic datasets
- 3. Compute t for each dataset, and report the spread of t for example with a histogram or a confidence interval

Exercise 9.2.3 (Comparing groups).

We are given data $x = [x_1, ..., x_m]$ which we believe is $N(\mu, \sigma^2)$ and further data $y = [y_1, ..., y_n]$ which we believe is $N(\mu + \delta, \sigma^2)$. Find a 95% confidence interval for $\hat{\delta}$.

The MLEs for μ , δ , σ are what you calculated in Example Sheet I question 4:

```
\hat{\mu} = \bar{x}
\hat{\delta} = \bar{y} - \bar{x}
\hat{\sigma} = \cdots
                                  \mathbf{x} = [4.3, 5.1, 6.1, 6.8, 7.4, 8.8, 9.9]
                                  y = [8.3, 8.5, 8.9]
                                  3 m,n = len(x), len(y)
                                  4 # 1. Define the readout statistic
                                  5 def t(x,y): return np.mean(y) - np.mean(x)
                                  7 # 2. To generate a synthetic dataset ...
                                  \hat{\mu}, \hat{\delta} = \text{np.mean}(\mathbf{x}), \text{np.mean}(\mathbf{y}) - \text{np.mean}(\mathbf{x})
                                      \hat{\sigma} = \text{np.sqrt}((\text{np.sum}((\mathbf{x}-\hat{\mu})**2 + \text{np.sum}((\mathbf{y}-\hat{\mu}-\hat{\delta})**2))/(\text{m+n}))
                                 10 def rxy_star():
                                            return (np.random.normal(loc=\hat{\mu}, scale=\hat{\sigma}, size=m),
                                                       np.random.normal(loc=\hat{\mu} + \hat{\delta}, scale=\hat{\sigma}, size=n))
                                 12
                                 13 # 3. Sample the readout statistic, and report its spread
                                 14 t_{-} = [t(*rx_star()) \text{ for } _in \text{ range}(10000)]
                                 15 lo, hi = np.quantile(t_{-}, [.025, .975])
                                 16 plt.hist(t_)
```



When you fit a model by maximizing Pr(data; params) use ALL the data and ALL the parameters.

When you resample, resample the FULL data.

GENERALIZATION, ACCORDING TO FREQUENTISTS

I trained my ML system on a dataset. How will it work on in-the-wild data?





- The job of a ML system is to report an output, e.g. a parameter or a prediction.
- I want to know what output my system might report on in-the-wild data.
- Let me consider an *ensemble* of systems, trained on many possible datasets. What range of outputs do they report?
- This range says how confident I can be about my system's output (assuming I've even got the probability model right).

My system reports a parameter estimate $\hat{\theta}$. What's the <u>frequency</u> of $\hat{\theta} \in [lo, hi]$ across the multiverse?

Office hours 1–1.30pm in the cafe area today

No lecture on Friday
Only Mon+Wed from now on