Syllabus for IB Data Science

Using a probability model to describe data

Models that depend on linear combinations of features

Parameter interpretation and identifiability

Fitting a model to the data

Maximum likelihood estimation

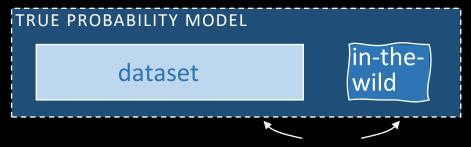
Fitting via least squares

Does the model generalize?



GENERALIZATION

I trained my model on a dataset. Will it work well on in-the-wild data?



We'll assume that dataset & in-the-wild satisfy the same Laws of Nature. (Otherwise it's impossible to say anything useful.)

Table 2: Results on HotpotQA distractor (dev). (+hyperlink) means usage of extra hyperlink data in Wikipedia. Models beginning with "—" are ablation studies without the corresponding design.

Model	Ans EM	Ans F_1	Sup EM	Sup F_1	Joint EM	Joint F_1
Baseline [53]	45.60	59.02	20.32	64.49	10.83	40.16
DecompRC [29]	55.20	69.63	N/A	N/A	N/A	N/A
QFE [30]	53.86	68.06	57.75	84.49	34.63	59.61
DFGN [36]	56.31	69.69	51.50	81.62	33.62	59.82
SAE [45]	60.36	73.58	56.93	84.63	38.81	64.96
SAE-large	66.92	79.62	61.53	86.86	45.36	71.45
HGN [14] (+hyperlink)	66.07	79.36	60.33	87.33	43.57	71.03
HGN-large (+hyperlink)	69.22	82.19	62.76	88.47	47.11	74.21
BERT (sliding window) variants						
BERT Plus	55.84	69.76	42.88	80.74	27.13	58.23
LQR-net + BERT	57.20	70.66	50.20	82.42	31.18	59.99
GRN + BERT	55.12	68.98	52.55	84.06	32.88	60.31
EPS + BERT	60.13	73.31	52.55	83.20	35.40	63.41
LQR-net 2 + BERT	60.20	73.78	56.21	84.09	36.56	63.68
P-BERT	61.18	74.16	51.38	82.76	35.42	(2.70
EPS + BERT(large)	63.29	76.36	58.25	85.60	41.30	
CogLTX	65.09	78.72	56.15	85.78	1	
- multi-step reasoning	62.00	75.39	51.74	83.10	1	
 rehearsal & decay 	61.44	74.99	7.74	47.37		
 train-test matching 	63.20	77.21	52.57	84.21	170	

Results. Table 2 shows that CogLTX outperforms most of previous method solutions on the leaderboard.⁴ These solutions basically follow the frame results from sliding windows by extra neural networks, leading to bounded to insufficient interaction across paragraphs.

Most ML papers don't state an inductive claim.

Perhaps the authors haven't thought hard enough to be able to state one?

Perhaps they prefer to leave you, the reader, to make the inference?

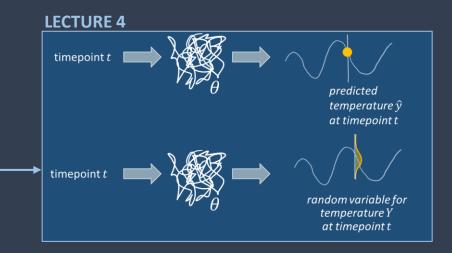
"All science is either physics or stamp-collecting."

Ernest Rutherford (1871-1937)

FITTING

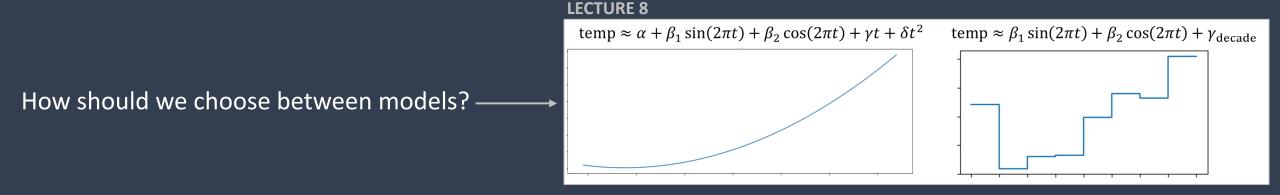
How do we measure how well a model "works"?

- prediction accuracy, or mean square error [MLRD]
- log likelihood is a more general metric -



So, let's fit our model's parameters by maximimum likelihood estimation. Then it'll surely work well on in-the-wild data!

FITTING AND MODEL CHOICE



Choosing between models can be seen as estimating a discrete parameter ...

1. Define a full model that has a "switch" parameter m specifying which submodel to use

```
\begin{array}{ll} \operatorname{def} \ \operatorname{ry}(\mathtt{m},\theta_1,\theta_2) \colon \\ & \text{if} \ \mathtt{m} == 1 \colon \\ & \operatorname{return} \ \operatorname{ry1}(\theta_1) \\ & \operatorname{else} \colon \# \ \mathtt{m} == 2 \\ & \operatorname{return} \ \operatorname{ry2}(\theta_2) \end{array} \qquad \operatorname{Pr}_Y(y;m,\theta_1,\theta_2) = \begin{cases} \operatorname{Pr}_Y^1(y;\theta_1) & \text{if} \ m = 1 \\ \operatorname{Pr}_Y^2(y;\theta_2) & \text{if} \ m = 2 \end{cases}
```

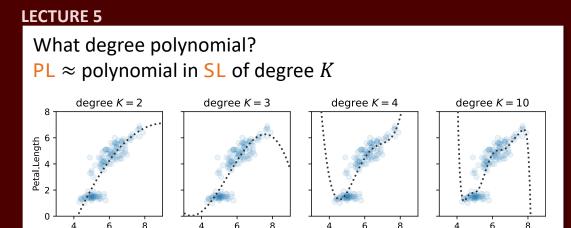
2. Then we can estimate it using maximum likelihood estimation: $\max_{m,\theta_1,\theta_2} \Pr(\text{data}; m, \theta_1, \theta_2)$

$$= \max \left\{ \max_{\theta_1} \Pr^1(\text{data}; \theta_1), \max_{\theta_2} \Pr^2(\text{data}; \theta_2) \right\}$$

The mle for m is the model with the highest likelihood, $\widehat{m} = \arg \max_{m} \Pr(\text{data}; m)$

MODEL COMPLEXITY

When we're choosing between models of different complexity, the more complex model will always score better.

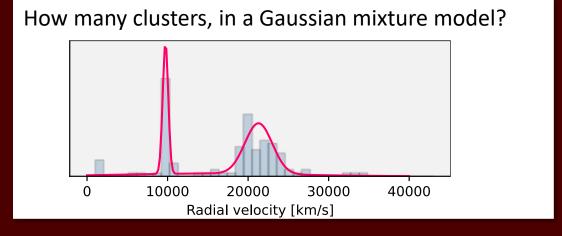


What model for the number of heads? $X \sim \text{Bin}(n, \frac{1}{2})$ • versus • $X \sim \text{Bin}(n, \theta)$ for some $\theta \in [0,1]$

LECTURE 8

Sepal.Length

LECTURE 3



OVERFITTING AND HOLDOUT EVALUATION

A too-complex model will typically fit the dataset well, but generalize poorly to in-the-wild data.

dataset

in-thewild

How can we detect and avoid too-complex models?

We can approximate in-the-wild performance by holdout performance.

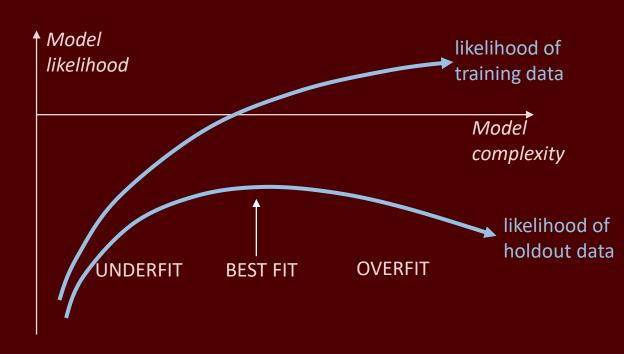
training

holdout in-thewild

This suggests we should simply choose the model that maximizes the holdout likelihood

With this approach, there's an inner loop where we fit a model to the training data, and an outer loop where we choose between fitted models.

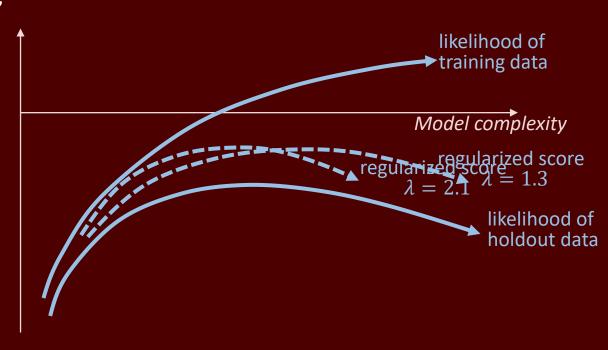
QUESTION. Can we do it all in a single loop?



REGULARIZATION

Some models have a natural continuous measure of complexity. In such cases, we can maximize the *regularized* score, which penalizes complex models: $\max_{\theta} \{\log \Pr(\text{data}; \theta) - \lambda \text{ complexity}(\theta)\}$

- Gradient descent will choose a tradeoff between fit and complexity.
 We're baking "too much complexity is bad" into the training itself.
- We control the relative importance by specifying $\lambda > 0$.
- (We'll still want to tune λ to maximize holdout likelihood, but this maximization is pretty forgiving.)



BAYESIANIST MODEL CHOICE

Suppose we want to choose between several models. We're uncertain which is the correct model.

BAYESIANIST EPISTEMOLOGY



Whenever there's an unknown parameter, you should express your uncertainty about it by treating it as a random variable.

- Q. What are we uncertain about? Depends on the problem
- Q. How do we represent unknowns?

 Answer: As random variables, with a prior
- Q. What do we report?

 Answer: The posterior distribution of the quantity of interest
- Q. How do we find this?

 Answer: Using Bayes's rule

Exercise 8.3.3 (Bayesian classification)

There are two types of expense claims, legitimate and fraudulent. The legitimate claim sizes are $\sim \operatorname{Exp}(\lambda_L)$ and the fraudulent ones are $\sim \operatorname{Exp}(\lambda_F)$ where $\lambda_L = 0.1$ and $\lambda_F = 0.02$. In my prior experience, 99% of claims I've seen are legitimate. A new claim comes in, for an amount £x. Is it likely to be fraudulent?

What are we uncertain about?

whether the new claim is fraudulent

How do we represent uncertainty?

Let
$$M = \begin{cases} \ell & \text{if the new claim is legitimate} \\ f & \text{if it's fraudulent} \end{cases}$$

What is my prior?

$$Pr_M(\ell) = 0.99$$
 and $Pr_M(f) = 0.01$

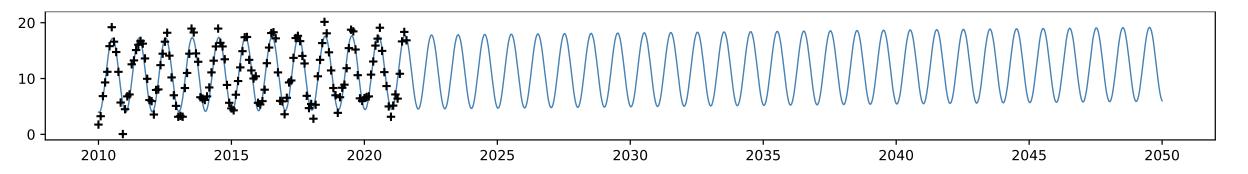
What is the posterior I want to report?

$$Pr_M(f \mid x)$$
 i.e. $\mathbb{P}(M = f \mid x)$

How should we choose between two models?

Modeller 1: Temp $\sim \alpha + \beta \sin(2\pi(t+\phi)) + \gamma(t-2000) + N(0,\sigma^2)$

Modeller 2: Temp $\sim \alpha' + \beta' \sin(2\pi(t + \phi')) + N(0, \sigma'^2)$



What are we uncertain about?

Which model is correct (and also all nine unknown parameters)

How do we represent uncertainty? With random variables.

Introduce a "switch" random variable M saying which model is correct, M=1 or M=2. Invent a prior for M, and for the other nine parameters. Pr(data | params) = Pr(temp₁,...,temp_n | M=m, $\alpha,\beta,\phi,\gamma,\sigma,\alpha',\beta',\phi',\sigma'$)

What do I want to report?

First find the posterior distribution of $(M,\alpha,...)$ given the data. Then report the marginal of M, $\mathbb{P}(M=1 \mid \text{data})$.

TIP: in multi-parameter problems, use Bayes's rule on *all* the unknowns simultaneously

our prior expresses our

preconceptions about

in-the-wild data

BAYESIANIST MODEL CHOICE

Suppose we want to choose between several models.

We're uncertain which is the correct model.

Introduce a "switch" parameter. Let the correct model be M, with prior $Pr_M(m)$.

$$Pr_M(m|data) = \kappa Pr_M(m) Pr(data|m)$$

We could simply report our posterior for M.

Or we could report a point estimate. The MAP estimator says to pick the m that maximizes the posterior likelihood:

$$\widehat{m} = \underset{m}{\operatorname{arg \, max}} \Pr_{M}(m|\operatorname{data})$$

$$= \underset{m}{\operatorname{arg \, max}} \left\{ \log \Pr(\operatorname{data}|m) + \log \Pr_{M}(m) \right\}$$

a regularizer, that penalizes models that are inconsistent with our prior beliefs

IB Data Science syllabus

Using a probability model to describe data

How well does my model fit the dataset?

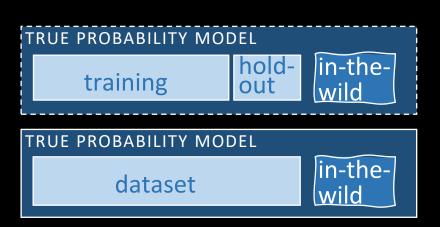
Note: model choice and parameter estimation are the same thing

How well does my model generalize to in-the-wild data?

[EMPIRICIST] We can use holdout evaluation to approximate in-the-wild performance, without making any assumptions (other than that there *is* a common true distribution).

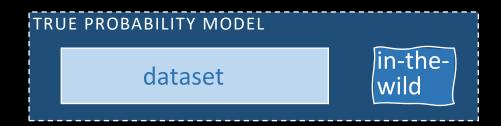
[BAYESIANIST] If we have a prior belief about the true distribution, we should use it. In this case, we don't need holdout evaluation.

[FREQUENTIST / HYPOTHESIS TESTING] ...

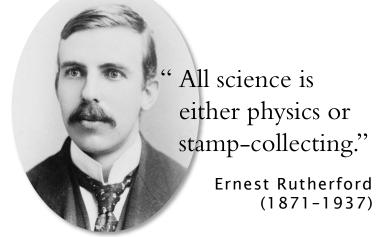


INDUCTION IS NOT DEDUCTION

I have data, and I've learned something from it. What can I say about the future? Philosophers call this *The Problem of Induction*.



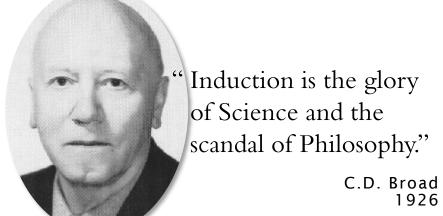
It is IMPOSSIBLE to mathematically *deduce* the true probability model from a dataset, so don't hanker after rigorous proofs for why one school of induction is better than another.





What's the chance the sun will rise tomorrow? $\frac{x+1}{n+2}$

Pierre-Simon Laplace (1749-1827)



IB Data Science syllabus — maths skills

Using a probability model to describe data

- familiarity with a range of models
- linear models; parameter interpretation

How well does my model fit the dataset?

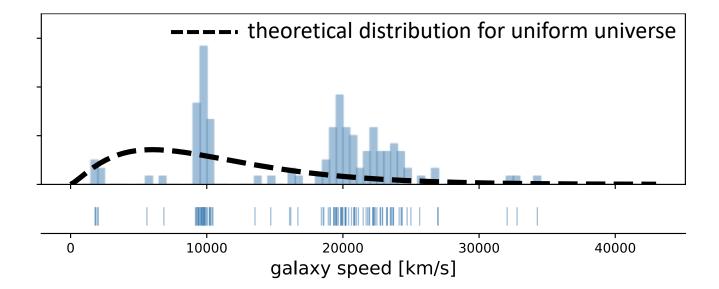
- maximum likelihood estimation
- numerical optimization

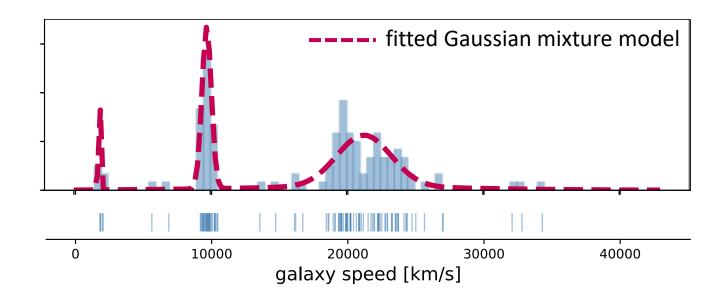
How well does my model generalize to in-the-wild data?

- Bayes's rule
- Monte Carlo
- bespoke distributions

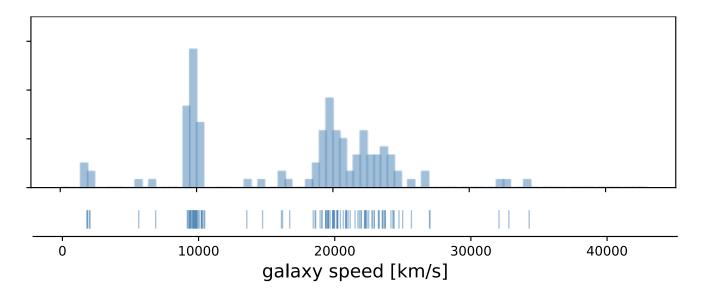
This chart shows the distribution of the speeds of 120 galaxies, from a survey of the Corona Borealis region.

Postman, Huchra, Geller (1986)



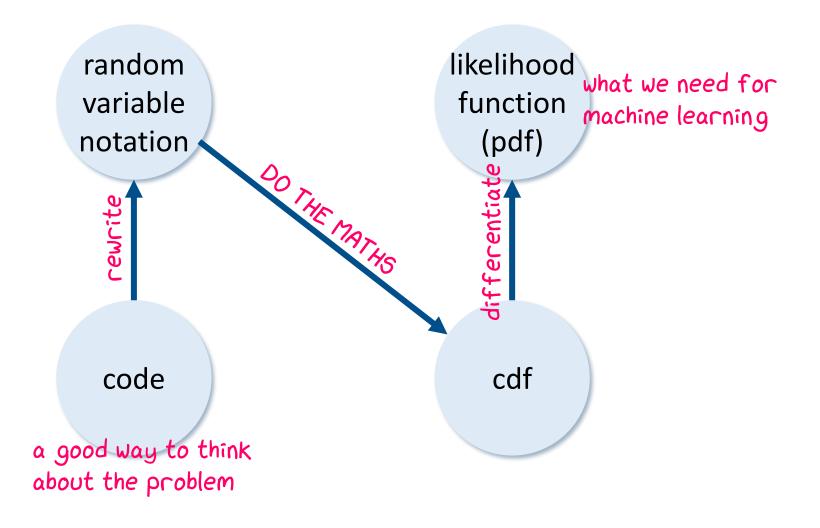


What's the best distribution we can find, to model this dataset?



Bespoke probability distributions part I: from code to likelihood (for continuous random variables)

There are four ways to specify a distribution:



Find the pdf of the random variable generated by this code:

$$x = - np.log(u) / \lambda$$



Step 1: random variable notation

Step 2: *X* is a continuous r.v., so find its cdf

Step 3: differentiate cdf to get pdf

Try to write our probability in terms of simple standard random variable (for which we can look up the cdf)

Break it down so that the random variables are on the left

(so we can use the textbook cdf)

Wikipedia: Uniform distribution

