§1.3 Maximum likelihood estimation

All of machine learning is based on a single idea:

- 1. Write out a probability model -
- 2. Fit the model from data

This is behind

- A-level statistics formulae
- our climate model
- ChatGPT training

i.e. estimate the parameters using Maximum Likelihood Estimation (MLE) - typically with unknown parameters

The likelihood is the probability of seeing the data that we actually some.

It depends on the parameters.

Let's simply pick the parameters that ^{§1} maximize the likelihood !



Likelihood of the observed data:

$$lik = IP(X = x)$$
$$= \binom{n}{x} P^{x} (I-P)^{n-x}$$

Parameter that maximizes it:

$$\frac{d}{dp} lik = \binom{n}{x} \left[x p^{x-1} (1-p)^{n-x} - (n-x) p^{x} (1-p)^{n-x-1} \right]$$

$$\frac{d}{dp} lik = 0 \qquad \Longrightarrow \qquad \hat{p} = \frac{x}{n}$$



There are standard numerical random variables that you should know:

DISCRETE RANDOM VARIABLES

| Binomial $X \sim Bin(n, p)$ | $\mathbb{P}(X = x) = \binom{n}{x} p^x (1-p)^{n-x}$ $x \in \{0, 1, \dots, n\}$ | For count data, e.g. number of heads in <i>n</i> coin tosses |
|---------------------------------------|---|--|
| Poisson $X \sim Pois(\lambda)$ | $\mathbb{P}(X = x) = \frac{\lambda^{x} e^{-\lambda x}}{x!}$ $x \in \{0, 1, \dots\}$ | For count data, e.g. number of buses passing a spot |
| Categorical $X \sim Cat([p_1,, p_k])$ | $\mathbb{P}(X = x) = p_x$ $x \in \{1, \dots, k\}$ | For picking one of a fixed number of choices |

CONTINUOUS RANDOM VARIABLES

| Uniform X~U[a, b] | $pdf(x) = \frac{1}{b-a}$ $x \in [a, b]$ | A uniformly-distributed floating point value |
|--|--|--|
| Normal / Gaussian $X \sim N(\mu, \sigma^2)$ | $pdf(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$ $x \in \mathbb{R}$ | For data about magnitudes, e.g. temperature or height |
| Pareto $X \sim Pareto(\alpha)$ | $pdf(x) = \alpha x^{-(\alpha+1)}$ $x \ge 1$ | For data about "cascade" magnitudes, e.g. forest fires |
| Exponential $X \sim Exp(\lambda)$ | $pdf(x) = \lambda e^{-\lambda x}$ $x > 0$ | For waiting times, e.g. time until next bus |
| Beta $X \sim Beta(a, b)$ | $pdf(x) \propto x^{a-1}(1-x)^{b-1}$ $x \in (0,1)$ | Arises in Bayesian inference |

There are standard numerical random variables that you should know:

Useful properties of the Normal distribution:

- If we rescale a Normal, we get a Normal
- If we add independent Normals, we get a Normal

$$a + b N(0,1) \sim a + N(0, b^2) - N(0, b^2)$$

for constants a and b
 $N(\mu, \sigma^2) + N(\nu, \rho^2) \sim N(\mu + \nu, \sigma^2 + \rho^2)$
assuming the two Normals are independent.

Normal / Gaussian

$$X \sim N(\mu, \sigma^2)$$
 pdf(x) = $\frac{1}{\sqrt{2\pi\sigma^2}}e^{-(x-\mu)^2/2\sigma^2}$ For data about magnitudes, e.g. temperature or height
 $x \in \mathbb{R}$



parts of this structure correspond

to reality.

The Nobel Prize in Chemistry 2024 Summary Laureates David Baker Demis Hassabis John M. Jumper Prize announcement

Press release

Popular information

Advanced information

Share this







Horizontal steam boiler, Augsburg machines, early 19th century



THÉORIE

ANALYTIQUE

DE LA CHALEUR,

PAR M. FOURIER.



A PARIS,

CHEZ FIRMIN DIDOT, PÈRE ET FILS, LIBRAIRES POUR LES MATHÉMATIQUES, L'ARCHITECTURE HYDRAULIQUE ET LA MARINE, RUE JACOB, Nº 24.

1822.

Two main results:

- any function of a variable, whether
 continuous or discontinuous, can be
 expanded in a series of sines of multiples
 of the variable → Fourier transform
- partial differential equation for conductive diffusion of heat

In "classic" data analysis, information propagates according to thermodynamics law

 \rightarrow tangent spaces to data manifold resembles Euclidean planes



Figure 2. A manifold \mathcal{M} and the vector space $T_{\mathcal{X}}\mathcal{M}$ (in this case $\cong \mathbb{R}^2$) tangent at the point \mathcal{X} , and a convenient side-cut. The velocity element, $\dot{\mathcal{X}} = \partial \mathcal{X} / \partial t$, does not belong to the manifold \mathcal{M} but to the tangent space $T_{\mathcal{X}}\mathcal{M}$.

J. Solà, J. Deray, D. Atchuthan, «A micro Lie theory for state estimation in robotics,», https://arxiv.org/abs/1812.01537, 2021



Exercise 1.3.1 (Coin tosses)

Suppose we take a biased coin, and tossed it n = 10 times, and observe x = 6 heads. Let's use the probability model

 $X \sim \operatorname{Binom}(n, p)$

where p is the probability of heads. Estimate p.

Log likelihood of the observed data:

$$lik = P(X = x) = {\binom{n}{x}} p^{x} (l - p)^{n-x}$$

log lik = log {\binom{n}{x}} + x \log p + (n-x) \log (l - p)

Parameter that maximizes it:

$$\frac{d}{dp} \log lik = \frac{x}{p} - \frac{n-x}{1-p}$$

$$\Rightarrow \hat{p} = \frac{x}{n}$$



Exercise 1.3.6 (Handling boundaries) We throw a k-sided dice, and get the answer x=10. Estimate k, using the probability model

$$\mathbb{P}(\operatorname{throw} x) = \frac{1}{k}, \quad x \in \{1, \dots, k\}$$

WHY??

lik

123

lik:
$$P(throw x) = \frac{1}{R}$$

 $k = 1$
 $SILLY$



SANITY CHECK Does our answer depend on the data? In the way we'd expect it to?



§1.3 Maximum likelihood estimation

All of machine learning is based on a single idea:

- 1. Write out a probability model -
- 2. Fit the model from data

This is behind

- A-level statistics formulae
- our climate model
- ChatGPT training

i.e. estimate the parameters using Maximum Likelihood Estimation (MLE) - typically with unknown parameters

The likelihood is the probability of seeing the data that we actually some.

> IP (data) if our model is a discrete rand.vour. pdf (data) if our model is a continuous rand.vour

⁵¹ If the data conjuts of many datapoints [X,,...,Xn], and our model says they're independent, lik (data) = lik (X1) × lik (X2) ×··· × lik (Xn)

It depends on the parameters.

Let's simply pick the parameters that maximize the likelihood!

Exercise 1.3.2 (Exponential sample)

Let the dataset be a list of real numbers, $x_1, ..., x_n$, all > 0. Use the probability model that says they're all independent $Exp(\lambda)$ random variables, where λ is unknown. Estimate λ .

Log likelihood of the observed data:

$$X \sim Exp(\lambda)$$

$$P(X = x_i) = 0$$

$$pdf(x_i) = \lambda e^{-\lambda x_i}$$

$$lik (dota) = lik (x_1) \times \cdots \times lik (x_n)$$

= $(\lambda e^{-\lambda \times i}) \times \cdots \times (\lambda e^{-\lambda \times i}) \times \cdots \times (\lambda e^{-\lambda \cdot \sum_{i=1}^{n} x_i})$
= $\lambda^n e^{-\lambda \cdot \sum_{i=1}^{n} x_i}$
loglik = $n \log \lambda - \lambda \cdot \sum_{i=1}^{n} x_i$

(CONTINUOUS RANDOM VARIABLES (real-valued) Exponential $pdf(x) = \lambda e^{-\lambda x}$ $X \sim Exp(\lambda)$ x > 0np.random.exponential(scale=1/ λ)

- Xxn

Parameter that maximizes it:

$$\frac{d}{d\lambda}\log(ik = \frac{n}{\lambda} - \sum_{i} x_i = 0 \implies \hat{\lambda} = \frac{n}{\sum_{i} x_i}.$$

Exercise 1.3.4 (Predictive models)

Consider a dataset of January temperatures, one record per year. Let t_i be the year for record i = 1, ..., n, and let y_i be the temperature. Using the probability model

$$Y_i \sim \text{Normal}(\alpha + \gamma t_i, \sigma^2)$$

estimate γ , the annual rate of temperature change.

Note: the question doesn't hell us the values of 02, 8, J, so we'll troat them as unknowns to be estimated.

$$\frac{\partial}{\partial \alpha} \frac{\log lik}{\partial \alpha} = 0$$

$$\frac{\partial}{\partial \delta} \frac{\log lik}{\partial \delta} = 0$$

$$\frac{\partial \log lik}{\partial \delta} = 0$$

$$\frac{\partial \log lik}{\partial \delta} = 0$$

§1.3 Maximum likelihood estimation



Exercise 1.3.4 (Predictive models)

Consider a dataset of January temperatures, one record per year. Let t_i be the year for record i = 1, ..., n, and let y_i be the temperature. Using the probability model

 $Y_i \sim \text{Normal}(\alpha + \gamma t_i, \sigma^2)$

estimate γ , the annual rate of temperature change.

What would happen if we just solved one equation, for the parameter we're interested in?

$$\frac{d}{l\gamma}\log lik = 0$$

We get the answer

$$\hat{\gamma} = \frac{\Sigma_i t_i (y_i - \alpha)}{\Sigma_i t_i^2}$$

SANITY CHECK Does our answer depend on unknown parameters?

Three views of a probability model



Temp:
$$N(\text{pred}; \sigma^2)$$

where $\text{pred}; = \alpha \sin(2\pi(f_i^+ \phi))$
 $+ c + \delta f_i^-$

$$\begin{split} \text{Temp}_i &\sim \alpha \sin \left(2\pi (t_i + \varphi) \right) + c + \gamma t_i + \text{Normal}(0, \sigma^2), \\ i &\in \{1, \dots, n\} \end{split}$$

The observed data is $[temp_1,...,temp_n]$. Find an expression for the log likelihood.

$$(ik (data) = lik (temp) \times \dots \times (ik (temp))$$

$$= \left(\frac{1}{12\pi\sigma^2} e^{-\frac{(temp)}{N} - pred}\right)^{2/2\sigma^2} \times \dots$$
Watch out for copy-paste-itis! We want the likelihood of seeing temp1, for the random variable Temp1~N(pred1, \sigma^2). Don't just paste in the formula from the random variable reference sheet,
$$\frac{1}{\sqrt{2\pi\sigma^2}}e^{-(x-\mu)^2/2\sigma^2}$$
log lik (data) = $-\frac{n}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum(temp) - predictioned$

RISOTTO CHALLENGE

• <u>Question</u>:

what is the major limitation of MLE in modern data analysis?

- Submit your answer (one per person) by <u>Nov. 1, 2024</u> sending an email to <u>am2920@cam.ac.uk</u> with subject "RISOTTO CHALLENGE"
- In the meanwhile, hints will be provided in the next classes
- No points earned for course assessment, but a valuable prize: THE risotto recipe from your dearest lecturer!
- Award ceremony on Class 11!

§1.4 Numerical optimization

All of machine learning is based on a single idea:

- 1. Write out a probability model
- 2. Fit the model from data —

This is behind

- A-level statistics formulae
- our climate model
- ChatGPT training

using maximum likelihood estimation with numerical optimization

(since the likelihood function is usually far too complex for exact optimization)

§1



+ There is no scipy.optimize.fmax

Exercise 1.4.2 (Constraints / softmax transformation) Find the maximum of

 $f(p_1, p_2, p_3) = 0.2 \log p_1 + 0.5 \log p_2 + 0.3 \log p_3$ over $p_1, p_2, p_3 \in (0,1)$ such that $p_1 + p_2 + p_3 = 1$.

Conning field:
inspead of finding wat ever
$$(p_1, p_2, p_3)$$
 such that $p_1 + p_2 + p_3 = 0$
we'll inspead find mat war $(s_1, s_2, s_3) \in \mathbb{R}^3$
and set $p_i = \frac{e^{s_i}}{e^{s_i} + e^{s_i} + e^{s_j}}$
This forces $p_i \in (0, 1)$, $p_1 + p_2 + p_3 = 1$
def f(p):
 $p_1, p_2, p_3 = p$
return 0.2*np.log(p_1) + 0.5*np.log(p_2) + 0.3*np.log(p_3)
def softmax(s):
 $p = np.exp(s)$
return $p / np.sum(p)$
 $\hat{s} = scipy.optimize.fmin(lambda s: -f(softmax(s)), [0,0,0])$
 $\hat{s} = softmax(\hat{s})$
Optimization terminated successfully. Current function value: 1.02965. Iterations: 63.
Function evaluations: 120

array([0.19999474, 0.49999912, 0.30000614])

1

6

8 9

10

de

ŝ

ŝ

0p



How does it work? Animations by Lili Jiang, <u>Towards Data Science</u>



GRADIENT DESCENT

Find the gradient of the function, and take a step in the direction of steepest descent





Visualizing the Loss Landscape of Neural Nets

Li, Xu, Taylor, Studer, Goldstein (2018)

https://arxiv.org/abs/ 1712.09913



Software 1.0 is code we write. Software 2.0 is code written by the optimization based on an evaluation criterion (such as "classify this training data correctly"). It is likely that any setting where the program is not obvious but one can repeatedly evaluate the performance of it (e.g. — did you classify some images correctly? do you win games of Go?) will be subject to this transition, because the optimization can find much better code than what a human can write.