

An Introduction to Learning

CRT in AI 2020

Alessio Benavoli

Senior Lecturer

Computer Science and Information Systems (CSIS)

University of Limerick

alessio.benavoli@ul.ie

alessiobenavoli.com

10th November 2020



What is (Machine) Learning?

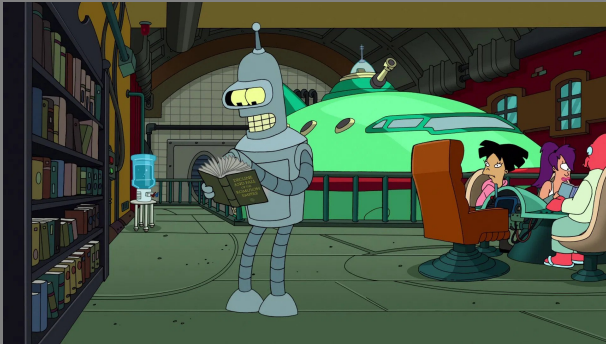


Table of Content

The number game

- Rule-based approach

- Ranked-hypothesis elimination

- Statistical learning

 - Likelihood

 - Prior

 - Posterior

 - Bayesian framework

 - Maximum Likelihood Estimation (MLE)

 - Maximum A-Posteriori (MAP)

How is this related to Machine Learning?

- Polynomial Regression

The number game

1. I tell you I am thinking of some simple arithmetical concept C .
2. I give you a series of **randomly** chosen positive examples

$$X = \{x_1, \dots, x_n\},$$

that is numbers that belong to the concept C .

3. Then I ask you whether any other numbers (**test cases**) y belong to the concept C .

For simplicity, we assume that all numbers are between 1 and 100, so **your** task is to compute whether $y \in C$ given X , for $y \in \{1, \dots, 100\}$; we call it the **generalization task**.¹

¹This first part of the talk is a summary from Ch. 1 and 2 of Joshua Tenenbaum's thesis, "A Bayesian Framework for Concept Learning".

Example

I am thinking about “Fibonacci numbers” (but you do not know it) and I give you three examples of the concept:

1, 5, 13

You have to guess the arithmetical concept C I am thinking about.

Note that many concepts are compatible with the example, for instance

odd numbers

The **generalization task** consists in giving examples of other number belonging to the concept, for instance

7

(this is a wrong generalization because 7 is not Fibonacci).

Let us start the game

I tell you 16 is positive example of the concept.

What other numbers do you think are positive examples?

Is 17? 4? 32? 87? 31?

Example

Is 17? 4? 32? 87? 31?

Maybe, it seems that 17 is more likely to be a positive example of the concept than 87, based on its **proximity** to 16.

Maybe 32 is more likely to be a positive example of the concept than 31, because it shares more **arithmetical** properties with 16.

Similarly, 4 seems to me more likely than 5 or 6.

Maybe 6 is more likely than 4, because it shares the same last digit as 16.

Example

More examples of the concept:

16, 8, 2, 64

Now your task should be easier!

Example

More examples of the concept:

16, 8, 2, 64

Now your task should be easier!

For instance, it seems that

- ▶ 32 and 4 become much more likely to be positive examples than 31 or 6;
- ▶ both 17 and 87 seem quite unlikely.

Example

The majority would agree that given the four **random** positive examples

16, 8, 2, 64

the concept program seems most likely to select the powers of two,

2, 4, 8, 16, 32, 64

and so on.

Surprising

It is quite amazing that we can conclude that, because given 100 numbers

$$\{1, 2, 3, \dots, 100\}$$

there are

$$2^{100}$$

possible subsets of numbers between 1 and 100.

Every time we give a positive example of the concept we cut in half the number of logically consistent subsets (that is those containing all the positive examples).

That means that there are still more than a billion billion subsets of numbers including 16, 8, 2, 64.

Example of these subsets

Example 16, 8, 2, 64:

- ▶ *all powers of two*
- ▶ *all even numbers*
- ▶ *all numbers less than 65*
- ▶ *all numbers less than 90*
- ▶ *all powers of two and also 25*
- ▶ *all powers of two except 32*
- ▶ ...

Despite this HUGE range of possibilities, we believe that we can identify numbers in the one subset that this program selects.

This belief comes after seeing just four random examples known to be in that set

- ▶ positive examples of the concept
- and no negative examples.

A first solution to the concept learning problem

A rule-based approach² can provide a basic solution to our concept learning problem.

In a rule-based system, we

1. consider various hypotheses about what the concept could be (a **reasonable** subset of the 2^{100} possible subsets)
2. eliminate hypotheses which are not consistent with the examples observed.

We denote the initial set of reasonable hypotheses by

$$\mathcal{H}$$

and we call it the **hypothesis space**!

Since the size of \mathcal{H} is less than 2^{100} , generalization from few positive examples becomes practical because most of the many ways to generalize are never even considered by the learner.

²The first AI expert systems were rule-based

Rule-based approach

Since the size of \mathcal{H} is much less than 2^{100} , generalization from few positive examples becomes practical because most of the possible ways to generalize are never even considered by the learner.

$\mathcal{H} = \{ \text{all powers of two, all powers of three, all powers of four, all even numbers, all odd numbers, all numbers less than 80, ...} \}$

After we see 16,

$\mathcal{H}_1 = \{ \text{all powers of two, } \cancel{\text{all powers of three}}, \cancel{\text{all powers of four}}, \text{all even numbers, } \cancel{\text{all odd numbers}}, \text{all numbers less than 80, ...} \}$

After we see 16, 8, 2, 64,

$\mathcal{H}_4 = \{ \text{all powers of two, } \cancel{\text{all powers of three}}, \cancel{\text{all powers of four}}, \text{all even numbers, } \cancel{\text{all odd numbers}}, \text{all numbers less than 80, ...} \}$

Rule-based approach

The observed examples 16, 8, 2, 64 are consistent with more than one a priori reasonable hypothesis and, therefore, it is quite unlikely that we can end with a single hypothesis after few examples!

Theoretically, this approach can only work if we have a small hypothesis space and a large enough set of examples to guarantee that only one hypothesis will survive elimination.

After we see 16, 8, 2, 64,

$\mathcal{H}_4 = \{ \text{all powers of two, } \text{all powers of three, } \text{all powers of four,}$
 $\text{all even numbers, } \text{all odd numbers, } \text{all numbers less than 80, } \dots \}$

What do we choose?

Ranked-hypothesis elimination

The rule-based approach is too simple. Maybe it can be improved if the initial hypothesis space are allowed to be soft, not hard.

That is, instead of just allowing or disallowing candidate hypotheses, we can rank all possible hypotheses using some a-priori idea of **reasonableness**, and the learner chooses the highest ranked hypothesis that is consistent with the examples.

After we see 16, 8, 2, 64 and given the hypothesis space:

$\mathcal{H} = \{ \text{all powers of two, all powers of three, all powers of four, all even numbers, all odd numbers, all numbers less than 80, } \dots \}$

under the ranked rule-based view, to justify the preference for *powers of two* over *even numbers*, we would have to assume that *powers of two* has a-priori higher rank than even numbers.

However there is a problem, how can we explain that after seeing

16, 8

and then

16, 8, 2, 64

in the second case, we are much more **sure** that *powers of two* is more likely than *even numbers*. The **ranked rule-based** cannot explain that.

Summarising

Both

- ▶ rule-based
- ▶ rank rule-based

are not good enough!.

We need a different way of learning:

- ▶ a model that is able to explain why some consistent hypotheses seem increasingly more likely than others as more examples are observed
- ▶ how generalization can be based on multiple hypotheses of (perhaps) varying plausibilities

The solution

1. how does the learner infer, from a small set of positive examples of a concept, which other objects are likely to fall under that concept?
2. why do we infer that, given the random yes examples of 16, 8, 2, 64, the program probably accepts all and only the powers of two?
3. First, why should the learner settle on one hypothesis, e.g. all powers of two, over others that are equally consistent with the observations 16, 8, 2, 64 and that seem equally or more natural a priori, such as all even numbers or all numbers less than 100?
4. Second, how should the learner generalize when no single hypothesis is clearly more compelling than all others, as after the one example 16?

Avoiding suspicious coincidences

Why do we choose

$h_1 = \text{powers of two}$ and not for instance $h_2 = \text{even numbers}$

after seeing $X = \{16, 8, 2, 64\}$, given that both hypotheses are consistent with the evidence?

The key idea is that we want to **avoid suspicious coincidences**.

If the true concept were *even numbers*, why we did not see any numbers that were not *powers of two*?

Avoiding suspicious coincidences

The fact that $X = \{16, 8, 2, 64\}$ is considered “suspicious” is because we are implicitly making the **strong sampling assumption**, that is that the examples were chosen randomly from the concept.

Under the **strong sampling assumption**, the probability of independently sampling n items (with replacement) from h is given by

$$\textbf{Likelihood: } p(X|h) = \left[\frac{1}{\text{size}(h)} \right]^n$$

This equation is called the **size principle**: it is a form of **Ockhams razor**, which says one should pick the simplest explanation that is consistent with the data.

Example

Let $X = \{16\}$, then

$$p(X|h = \text{powers of two}) = \frac{1}{6}$$

because there are only 6 powers of 2 less than 100.

This is more likely than:

$$p(X|h = \text{even numbers}) = \frac{1}{50}$$

and much more likely than inconsistent concepts

$$p(X|h = \text{odd numbers}) = 0.$$

Example

After seeing $X = \{16, 8, 2, 64\}$,

$$p(X|h = \text{powers of two}) = \frac{1}{6^4} = 7.7 \cdot 10^{-4}$$

This is much more likely than:

$$p(X|h = \text{even numbers}) = \frac{1}{50^4} = 1.6 \cdot 10^{-7}$$

This is a likelihood ratio of almost 5000 : 1 in favor of *powers of two*. This confirms mathematically our earlier intuition that $X = \{16, 8, 2, 64\}$ would be a very suspicious coincidence if generated by *even numbers*.

What is the least suspicious hypothesis?

The is the hypothesis that has the maximum likelihood given $X = \{16, 8, 2, 64\}$:

$$\max_h p(X|h) = \frac{1}{\text{size}(h)^4}$$

Note that the most likely hypothesis is not *powers of two*! For instance, the rather unreasonable hypothesis, *powers of two except 32* has

$$p(X|h = \text{powers of two except } 32) = \frac{1}{5^4} = 1.6 \cdot 10^{-3}$$

It has higher likelihood because it does not need to explain the (small) coincidence that we did not see 32.

To rule out such unreasonable concepts, we need to use some prior knowledge.

Prior

After seeing $X = \{16, 8, 2, 64\}$, Why do we (humans) choose

$h_1 = \text{powers of two}$

and not, for instance,

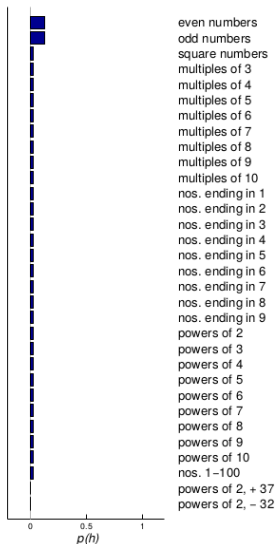
$h_2 = \text{powers of two except } 32$

After all, h_2 has higher likelihood!

However, h_2 is much less likely than h_1 a-priori, because it is **conceptually an unnatural hypothesis**.

It is the combination of the likelihood and the prior knowledge that can explain why we (humans) select h_1 .

Prior

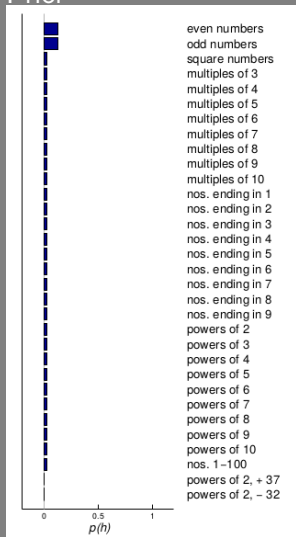


This is my prior, it puts less weight on unnatural concepts such as *powers of two except 32*, and more weight on very simple concepts like *even numbers*.

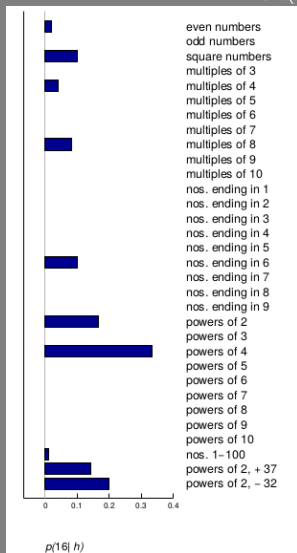
You may have a different prior. In general, the prior encodes the background knowledge on the “concept” we aim to learn.

Prior and Likelihood

Prior

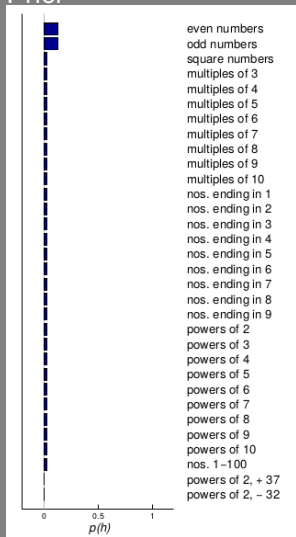


Likelihood ($p(16|h) = \frac{1}{\text{size}(h)}$)

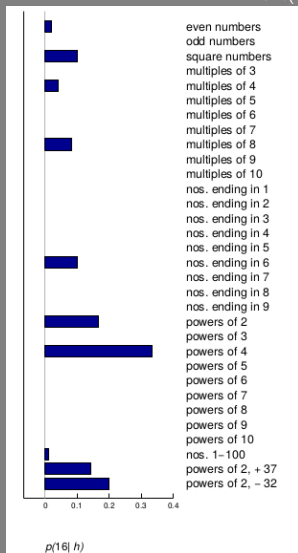


How do we combine them?

Prior



Likelihood ($p(16|h) = \frac{1}{\text{size}(h)}$)



How do we combine them?

We have defined

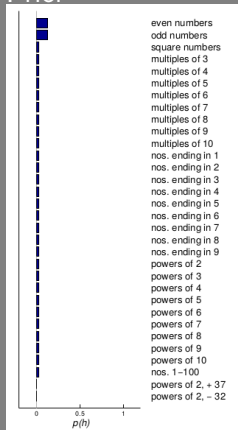
- ▶ $p(h)$ is the prior probability on the hypotheses, it tells us how probable we think it is that h is the true hypothesis before we have observed any examples.
- ▶ $p(X|h)$ is likelihood, it tells us the probability that we would observe the examples X if h were the true hypothesis.
- ▶ we want $p(h|X)$, the probability that the true hypothesis is h given we have observed X . This is called posterior and measures our belief in h after we observed the examples in X .

We want to compute:

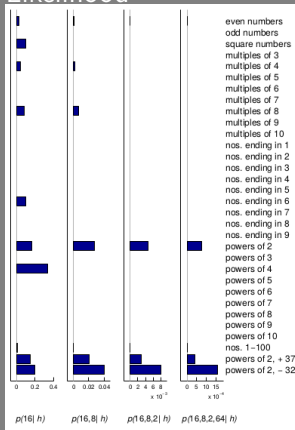
$$p(h|X) = \frac{p(X|h)p(h)}{p(X)} = \frac{p(X|h)p(h)}{\sum_{h'} p(X|h')p(h')}$$

that is Bayes' rule.

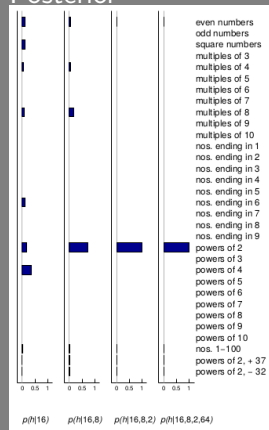
Prior



Likelihood



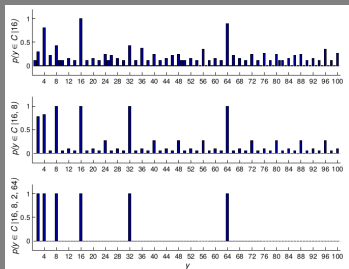
Posterior



Problem of generalization

In the Bayesian setting, the problem of generalization reduces to computing $p(y \in C|X)$, that is the probability that a new number y belongs to concept C , given the set X of previously observed. (Bayesian) Model averaging:³

$$p(y \in C|X) = \sum_h p(y \in C|h)p(h|X)$$



³ $p(y \in C|h) = 1$ if $y \in h$ and zero otherwise.

Summarising the Bayesian framework

1. A constrained hypothesis space \mathcal{H} . Otherwise, it is impossible to generalize from a finite data set, because any hypothesis consistent with the evidence is possible.
2. An informative prior, that ranks members of the hypothesis space.
3. The size principle, which is the likelihood function of a strong sampling model.
4. Hypothesis averaging, i.e., integrating out h when making predictions:

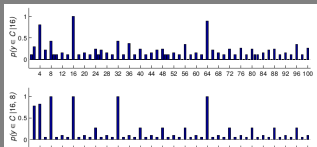
$$p(y \in C|X) = \sum_h p(y \in C|h)p(h|X)$$

Maximum Likelihood Estimation (MLE)

It has only ingredients 1 and 3.

- ▶ 1. A constrained hypothesis space \mathcal{H} . Otherwise, it is impossible to generalize from a finite data set, because any hypothesis consistent with the evidence is possible.
- ▶ 3. The size principle, which is the likelihood function of a strong sampling model.

Issues: given the example 16, the MLE hypothesis is *all powers of four*, so only 4 and 64 receive a nonzero probability of generalization, while given two examples $\{16, 8\}$, all 4, 8, 16, 64 receive a nonzero probability of generalization. MLE can make very sharp decisions after only one (few) observation (**overfitting**), while for the Bayesian method this is never the case:



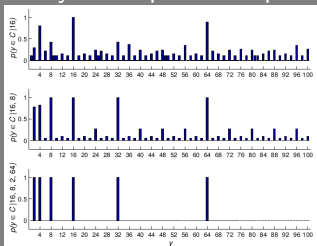
Maximum A-Posteriori (MAP)

It has ingredients 1, 2 and 3.

- ▶ 2. An informative prior, that ranks members of the hypothesis space.

It returns the hypothesis that has the highest posterior probability:
 $\max_h p(h|X)$.

Issues: without averaging, MAP does not account for the uncertainty, that is the fact that there may be alternative hypotheses that have very close posterior probability.



How is this related to Machine Learning?

The vast majority of the algorithms (for regression, classification and clustering) you have used (will use) are based on

- ▶ MLE (minimum loss function)
- ▶ MAP (regularisation)

Scikit-learn the most used machine learning library in Python:
Linear regression (MLE), Logistic regression (MAP), Neural Network (MAP)...

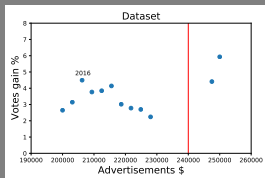
Keras is the most used deep learning framework (it uses MLE by default).

Issue: MLE leads to overfitting and MLE/MAP do not consider the uncertainty.

Parametric regression: polynomial model

2020 USA presidential election in Pennsylvania

1. In 2016 election, the Republican party spent about \$210,000 in advertisements, resulting in a 6% votes gain with respect to their baseline.
2. In 2020, the Republican party spent about \$240,000, resulting only in a 2.5% votes gain.⁴



Can we explain that based on historical data from past elections?

⁴This problem and dataset is not real, just an artificial example.

Is it a fraud?

We are going to use ML to answer this question.

Ingredients:

1. hypothesis space;
2. likelihood (suspicious coincidences);
3. prior on hypotheses (more reasonable/natural);
4. posterior;
5. averaging.

Ingredient 1: hypothesis space

Our goal is to find a relationship that explains how y (votes gain) depends on x (advertisements):

$$y = f(x)$$

We need to define a hypothesis space, for instance:

$$\mathcal{H} = \{f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_8 x^8, \beta_i \in \mathbb{R}\}$$

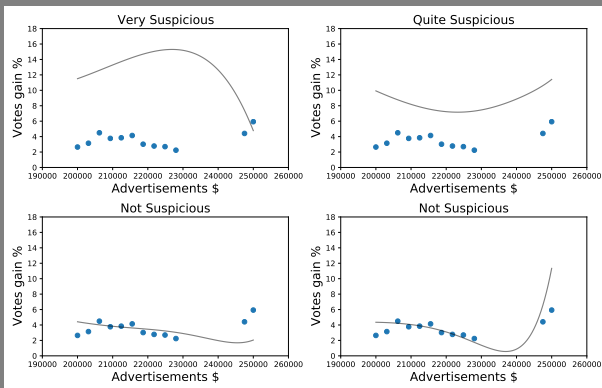
that is we assume that f is a polynomial function of degree 8. This is a large model class, for instance includes

1. linear relationship $f(x) = \beta_0 + \beta_1 x$
2. quadratic relationship $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$
3. ...

Ingredient 2: likelihood

Avoiding suspicious coincidences.

Likelihood $p(data|h) = p(data|\beta_0, \dots, \beta_8)$.



How do we define suspicious?

Closer to data means less suspicious:

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i - \beta_2 x_i^2 - \cdots - \beta_8 x_i^8)^2$$

this is called **squared error** loss function.

It is equivalent to assume a likelihood:

$$p(\text{data} | \beta_0, \dots, \beta_8) = \prod_{i=1}^N N(y_i; \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_8 x_i^8, \sigma)$$

which corresponds to the model

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_8 x_i^8 + \text{noise}$$

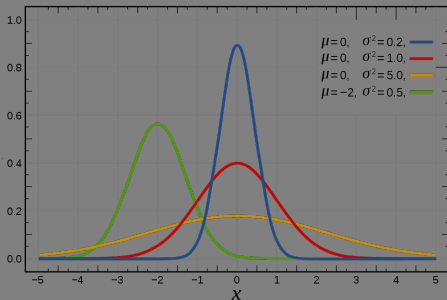
where noise is Gaussian with zero mean and standard deviation σ .

Gaussian PDF

Gaussian (or normal) distribution has PDF of a continuous variable x that can take values in the possibility space $\Omega = \mathbb{R}$ (a real number).

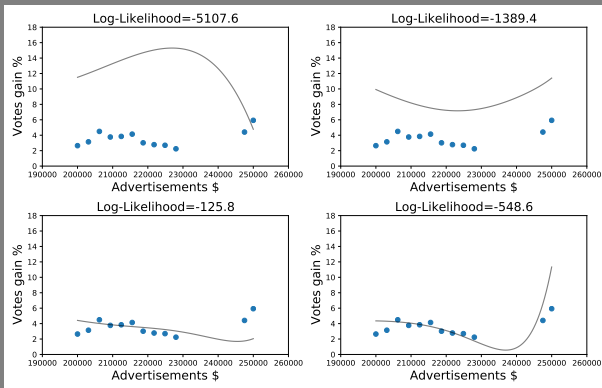
It is denoted as $N(x; \mu, \sigma)$ and defined as

$$p(x) = N(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$



How do we define suspicious?

We can use the likelihood!



Maximum Likelihood Estimator (MLE)

We determine the least suspicious hypothesis:

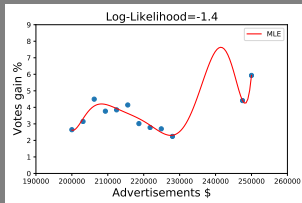
$$\min_{\beta_0, \dots, \beta_8} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i - \beta_2 x_i^2 - \dots - \beta_8 x_i^8)^2$$

which is equivalent to maximising the likelihood::

$$\max_{\beta_0, \dots, \beta_8} \prod_{i=1}^N N(y_i; \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_8 x_i^8, \sigma)$$

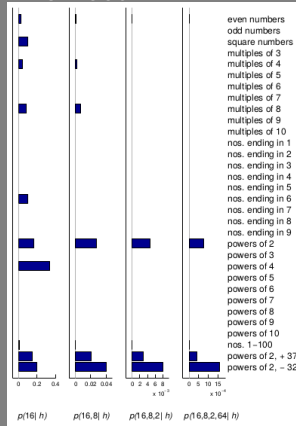
Solution:

$$[\beta_0, \dots, \beta_8] = [2.62, -4.93, 485.88, -4820.95, 21456.29, -51632.73, 68795.28, -47378.28, 13102.74]$$



MLE can overfit!

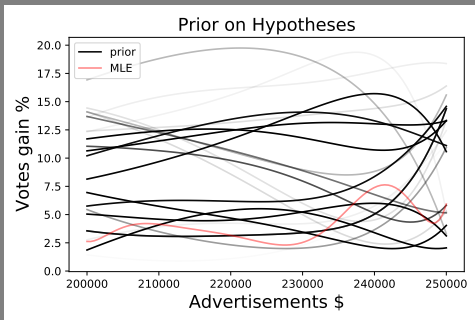
Likelihood



Think about the Number game, best MLE hypothesis was “power 2 excluded 32”.

Ingredient 3: Prior on hypotheses

$$p(h) = p(\beta_0, \dots, \beta_8) \quad \left(= N(\beta_i, 0, 10) \right)$$

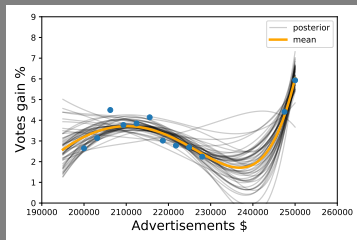


MLE polynomial looks quite unnatural (it has indeed very very small prior probability):

$$[\beta_0, \dots, \beta_8] = [2.62, -4.93, 485.88, -4820.95, 21456.29, -51632.73, 68795.28, -47378.28, 13102.74]$$

Ingredient 4: posterior

$$\underbrace{p(\beta_0, \dots, \beta_8, \sigma | \text{data})}_{\text{posterior}} = \frac{\underbrace{p(\text{data} | \beta_0, \dots, \beta_8, \sigma)}_{\text{likelihood}} \underbrace{p(\beta_0, \dots, \beta_8, \sigma)}_{\text{prior}}}{\underbrace{p(\text{data})}_{\text{evidence}}}$$



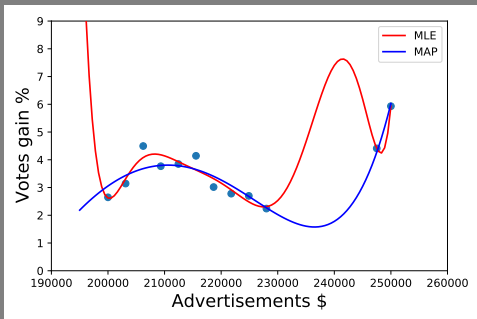
MAP (regularisation)

Instead of MLE we can compute the Maximum A-posterior hypothesis:

$$\max_{\beta_0, \dots, \beta_8} p(\beta_0, \dots, \beta_8, \sigma | \text{data})$$

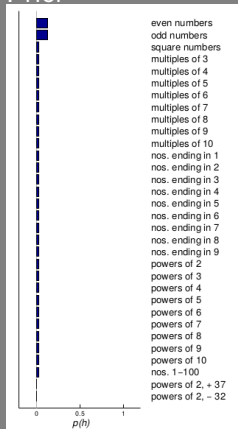
MAP polynomial:

$$[\beta_0, \dots, \beta_8] = [3.05, 7.1, -16.31, -3.34, 4.43, 6.21, 4.74, 1.81, -1.68]$$

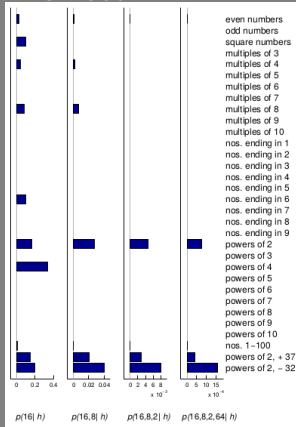


We need last ingredient to solve the last issue.

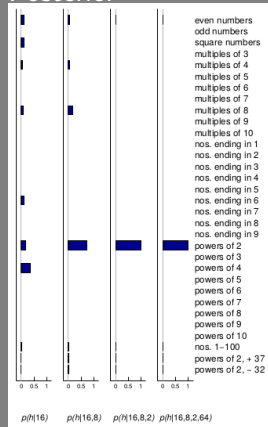
Prior



Likelihood

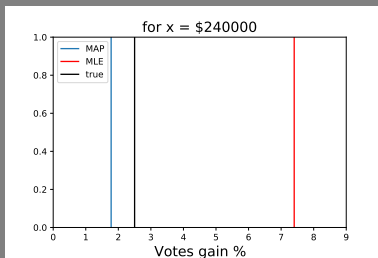


Posterior



Prediction

We are interested in predicting the “Votes gain” for $x = 240,000$:



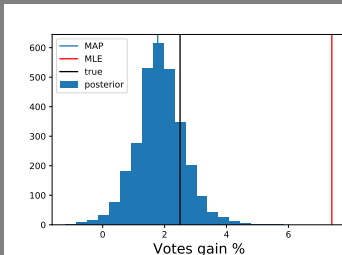
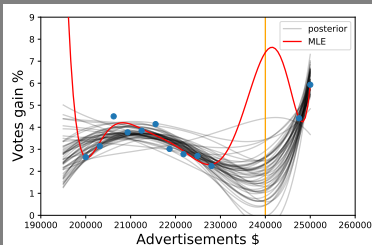
MLE estimate is **not reasonable**. The MAP estimate, that combines likelihood and prior, is closer to the true value (2.5%).

This explains why in Machine Learning, regularisation usually leads to more accurate prediction performance.

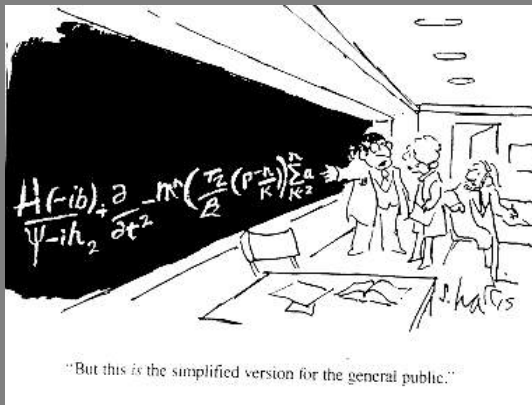
Ingredient 5: averaging

This does not answer our question, is the obtained 2.5% strange?
(Is it a fraud?)

We need can compute the distribution and see if 2.5% was strange. This shows us the uncertainty



How do we get those lines?



It is Bayesian Inference and as you learned we need to solve integrals. In this case, we can do it in closed form, but more in general we can use Stan or PyMC3 or any other probabilistic programming language.

How do we get those lines?

```
import pymc3 as pm
with pm.Model() as model:
    #prior
    theta_b = pm.Normal('theta', mu=0, sd=10, shape=(1+8,))
    sigma = pm.HalfCauchy('sigma', 1)
    #linear model
    mu = pm.Deterministic('mu', pm.math.matrix_dot(X8_tr,theta_b))
    #likelihood
    y_obs = pm.Normal('y',mu=mu, sd=sigma, observed=y_tr)
    posterior = pm.sample(3000,chains=1,tune=1000)
```

Conclusions

What did you learn?

the five ingredients for learning

How do you use what you learned?

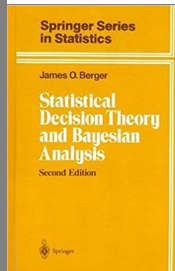
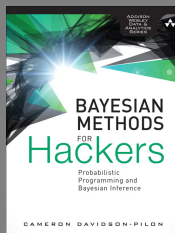
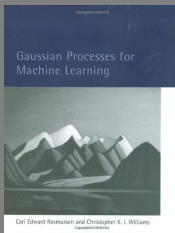
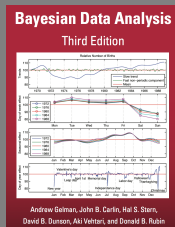
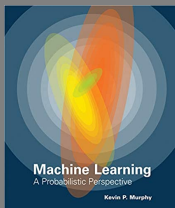
You can start this week: when you learn/use a new machine learning algorithm (NN, logistic regression, SVM...) you should always ask yourself “what is Ingredient 1?”, “What is Ingredient 2 (likelihood/loss function)?” etc.

There are two main approaches to Machine Learning:

- ▶ standard Machine Learning (MLE, MAP);
- ▶ if you are **bold** enough, you can decide to learn/use probabilistic Machine Learning (all five ingredients).

My favourite books

but also the books you should read



References

- ▶ The number game is from Joshua Tenenbaum's thesis, "A Bayesian Framework for Concept Learning", MIT, 1999.
- ▶ PyMC3 <https://docs.pymc.io/>