

Sampling Bias Correction for Supervised Machine Learning

A Bayesian Inference Approach with Practical Applications

Max Sklar

November 18, 2021

Abstract

Given a supervised machine learning problem where the training set has been subject to a known sampling bias, we ask if we can still train a model to fit the original dataset. This is achieved by introducing a sampling function to the Bayesian inference formula, and deriving an altered posterior distribution for the task. We then apply this to the common case of binary logistic regression, and discuss scenarios where a dataset might be subject to intentional sample bias such as label imbalance.

1 Introduction

A central task in computer science is developing *algorithms*. Each algorithm is a series of precise instructions to compute any given *function* $f : X \rightarrow Y$.

In *supervised machine learning*, a more general algorithm *learns* to compute a specific function f through example rather than through direct instructions from the programmer. These examples are given in the form of *instances* of input-output pairs $(x, y) \in X \times Y$ known as the *training set*. The training set allows the machine to learn the *concept* of f by producing an approximation of it known as a *model*. These models are often tuned on *validation sets* and evaluated on a *test set* of instances, and finally deployed on real world data to make predictions.

Machine learning is most reliable when the training set approximates these real world conditions. Models can in fact perform well on data that differ in composition from the training set, and engineers must sometimes rely on this fact. But that assumption is far from certain to hold true.

Occasionally an accurate training set is not available either by design or circumstance, and all that remains is a biased sample. **Fortunately, when the sampling function is known, there are methods to learn the concept from the remaining evidence available, and to design a system that builds a model for the original dataset.**

One method becomes apparent when analyzing the problem through the *Bayesian framework* for supervised machine learning reviewed in section 3. Key to this framework is the *posterior probability distribution* over potential models, which allows a search algorithm to discriminate between them by their projected performance. When a bias sample is provided, the posterior distribution formula can be properly corrected as detailed in sections 4 and 5 to counteract this bias.

2 Motivation

While big data has unlocked incredible new applications of machine learning, the accompanied processing comes at a heavy cost in terms of time, money, and energy. Many algorithms will perform just as well or almost as well had they been trained on a fraction of the data available. Therefore, removing data from the training set is a legitimate design decision.

One way to slim down the training set is through *uniform random sampling*, where each instance in the original dataset has an equal probability of inclusion. This process ignores the fact that some training examples are more valuable than others, particularly when there is a classification imbalance.

2.1 Example: Theoretical Image Recognition

Suppose that the goal is to train an image recognition algorithm to find images of lions. Furthermore, the training data contains 2000 images of lions, and 400,000 images that are **not** lions. The following 3 sampling schemes are proposed:

- (A) Train on all of the data (2000 Lions, 400000 Non-Lions)
- (B) Randomly sample 25 percent of the data (500 Lions, 100000 Non-Lions)
- (C) randomly sample a quarter of **ONLY** the non-lion images (2000 Lions, 100000 Non-Lions)

These 3 options each come with tradeoffs.

Option A uses all of available data and takes the most time and resources to train. In theory, the model produced will be at least as good as any of the others because the additional data can only refine our search for the best model. In practice this might not be the case because the search algorithm is imperfect. For example, if a *mini-batch* system is used to find gradients and search for the best model, it might not converge as easily with an imbalanced dataset.

Option B has all the imbalance of A, but with fewer instances. The learning algorithm will run more quickly, but the results are likely to be worse. Any classifier must rely on lion photos to learn lions, and those have become quite scarce! A reduction in lion photos from 2000 to 500 is likely to make a significant negative impact. Because of the imbalance, each lion photo is much more valuable to our success than a non-lion photo.

Under option C, the data size is almost as small as B, but now all 2000 lion photos remain. The likely outcome is that option C performs far better than option B. There's a good chance that it performs just as well as option A, but now the resources to run it will be as low as in option B.

One concern is that if option C does not correct for the sampling bias, it will use learn that Lions are more common than they really are in the underlying dataset, and will therefore tend to overpredict lions. Hopefully the underlying image recognition algorithm will derive the core visual features of a lion and will not have that issue even without correction, but such a correction is desirable.

2.2 Example: Probabilistic Event Detection

While bias correction in the image recognition example is desirable, in other cases it is crucial!

An iteration of Foursquare’s attribution model[13] provides a clear example. The Foursquare attribution product measures the ability of an advertisement to drive consumers to physical locations such as a retail chain. In order estimate causality, the Foursquare data pipeline must first learn the base probability that any given individual will visit that chain on a given day.

Foursquare’s data set has many examples of visits, but the amount of people who “did not visit” on any given day is higher by several orders of magnitude. Therefore, it makes sense to downsample those non-visits. Because a precise probability is necessary for an ad measurement product, that sampling must be accounted for appropriately. This was ultimately achieved through the bias correction formula for logistic regression in section 6.

2.3 Similar and Adjacent Work

The literature on sampling and rare events is vast, and no comprehensive review will be given here. Instead we refer to some selected work in this field which either inspired or was helpful in the course of this research.

For general cataloging of situations with imbalanced classification and techniques for it’s mitigation, see the work of Maalouf and Trafalis[9]. The readjustment formulas in the case of logistic regression can be found in Maalouf and Siddiqi[8]. For an in-depth discussion on rare events in logistic regression, the problems associated with it, and the mathematics of parameter estimation, see King and Zeng[6].

3 Supervised Machine Learning

The following is a typical setup for supervised machine learning using a Bayesian framework, a variation of what can be found in general treatments of probabilistic inference[10][4][2]. The terminology and variable names will be reused in subsequent sections.

3.1 Given Parameters

Let X be the input space and Y be the output space. The goal is to predict a *label* $y \in Y$ from a corresponding input $x \in X$, or to learn the function $f : X \rightarrow Y$.

The training dataset $D \in (X \times Y)^N$ consists of N examples of input-output pairs $(x, y) \in (X \times Y)$. D may be generated by an oracle which produces an arbitrary number of examples, or it might be a small and limited collection. In any case, inference will be based off of just these N examples.

Instances are labelled with a subscript (x_n, y_n) , where $n \in \{0, 1, 2, \dots, N - 1\}$, or $n \in N$ in ordinal notation.

3.2 The Hypothesis Space

Let H be the *hypothesis space* whose members $h \in H$ each encode a potential solution to the prediction problem. We assume that one of these solutions is correct. This temporary assumption makes Bayesian inference possible.

Each $h \in H$ is considered a *hypothesis* for the correct function f . They are also called *predictors*, *models*, and *solutions*. We will use the term predictor to emphasize its ultimate purpose.

In some setups, the predictors $h \in H$ are *directly predicting* $y \in Y$ and are represented by a function from X to Y . Here instead the predictors will return a *probability distribution* over Y , which means that $h \in H$ is now *probabilistically predicting* $y \in Y$. While it is possible to generalize to all *probability measures* over Y , we keep several cases in mind.

1. Y is finite. This is a problem of *classification*. Each predictor takes an input $x \in X$ and returns a probability for each $y \in Y$ that sums to one.
2. Y is discrete but infinite. The predictor still assigns a number to each potential output, but now the infinite sum adds to one.
3. Y is continuous. The predictor returns a *probability distribution function* (PDF) over Y , which integrates to one.

In each of these cases, the model assigns a number to each $y \in Y$. The case of discrete Y it is a finite probability and in the case of continuous Y it is the value of the PDF. If these numbers are normalized, the discrete probabilities will add to 1 and the continuous probabilities will integrate to 1.

These values may not be normalized at first. Instead we identify each predictor h with a *relative probability function* $f_h : X \times Y \rightarrow \mathbb{R}_{\geq 0}$ which assigns a non-negative real number to each input-output pair. Define the normalized probability function $\hat{f}_h : X \times Y \rightarrow \mathbb{R}_{\geq 0}$ as follows for both discrete and continuous Y .

$$\hat{f}_h(x, y) = \frac{f_h(x, y)}{\sum_{y' \in Y} f_h(x, y')} \quad \hat{f}_h(x, y) = \frac{f_h(x, y)}{\int_{y' \in Y} f_h(x, y')} \quad (1)$$

3.3 Deriving the Posterior Distribution from Bayes Rule

The *prior distribution* over H is a probability distribution representing the initial belief over which predictor is correct. Typically, this will involve defining an *uninformative prior* which encodes absence of knowledge of the problem. It is also common to incorporate *Occam's Razor* which penalizes more complex predictors, or to use a prior that is mathematically convenient. We will use $\mathbf{P}(h)$ as the *prior probability* of predictor $h \in H$. This may denote a discrete probability or a value in a PDF in the continuous case.

Fortunately, there is no need for $\mathbf{P}(h)$ to be normalized. It is enough for it to represent another relative probability function on h which at least preserves the ratio of the probabilities between two predictors. This includes the possibility of an *improper probability distribution function*, for example one that assigns a score of 1 to each real number if $H = \mathbb{R}$ even though such a PDF cannot be normalized. More importantly, it allows for distributions whose normalization is not tractable.

$\mathbf{P}(D|h)$ is the *likelihood function* which denotes the probability of receiving the entire training set D under a given predictor h . This is equal to the product of the probabilities of receiving each label y_n independently.

$$\mathbf{P}(D|h) = \prod_{n \in N} \hat{f}_h(x_n, y_n) \quad (2)$$

$\mathbf{P}(h|D)$ is the *posterior distribution* over H . This probability distribution represents the probability of each predictor being correct **after** the data has been taken into account.

Bayes rule for discrete H finds the posterior probability of individual predictor h , and in the continuous case it produces PDF values.

$$\mathbf{P}(h|D) = \frac{\mathbf{P}(D|h)\mathbf{P}(h)}{\sum_{h \in H} \mathbf{P}(D|h)\mathbf{P}(h)} \quad \mathbf{P}(h|D) = \frac{\mathbf{P}(D|h)\mathbf{P}(h)}{\int_{h \in H} \mathbf{P}(D|h)\mathbf{P}(h)} \quad (3)$$

Because most learning algorithms only require unnormalized values for $\mathbf{P}(h|D)$, we rewrite the equality as a proportionality statement and remove the denominator. This form also allows the prior $\mathbf{P}(h)$ to be unnormalized.

$$\mathbf{P}(h|D) \propto \mathbf{P}(D|h)\mathbf{P}(h) \quad (4)$$

The likelihood of the dataset under predictor h , or $\mathbf{P}(D|h)$, is equal to product of the absolute probability function of each instance. We rewrite the formula for the relative posterior distribution as

$$\mathbf{P}(h|D) \propto \left(\prod_{n \in N} \hat{f}_h(x_n, y_n) \right) \mathbf{P}(h). \quad (5)$$

3.4 Selecting and Sampling Predictors

The final learning task is either to *select* or *sample* a predictor that best explains the data. This process involves a search of H and is often identified as the search algorithm or learning algorithm.

In the case of selection, the goal is to identify a single optimal predictor. A good example is the *maximum a posteriori* (MAP) estimate which seeks to find the predictor with the highest posterior probability. Another example is the *maximum likelihood estimate* (MLE) which finds the predictor that assigns the highest likelihood to the dataset.

Under sampling, a predictor is randomly pulled from a distribution approximating the posterior distribution. Unlike selection, this method accounts for the uncertainty inherent in the result. If several models are sampled we can see how much variety there is in the possible solutions still consistent with the data.

In general, the hypothesis space H is complex and a variety of methods for selection and sampling have been developed. These techniques include hill climbing and gradient descent for selection, and Markov

Chain Monte Carlo for sampling. A good example of the latter is the *No U-Turn Sampler*[5] popular in the PyMC3[11] probabilistic programming package for python. This author has also deployed the Newton Raphson method for gradient descent to calculate the MLE.[12]

For most of these algorithms, it is more convenient to work with a *negative log-likelihood loss* function than the posterior distribution directly.¹ This is because the negative-log likelihood turns products into sums, and produces values that are within a reasonable order of magnitude for computation. We get this by applying $-\ln(\dots)$ to the right hand side of equation 5.

$$L(h) = - \sum_{n \in N} \ln(\hat{f}_h(x_n, y_n)) - \ln(\mathbf{P}(h))$$

Assume each hypothesis comes with its own negative log-likelihood loss function l_h where $\hat{f}_h(x_n, y_n) \propto e^{-l_h(x_n, y_n)}$ and the prior $\mathbf{P}(h)$ can be reduced to a *regularization function* $\mathbf{r}(h)$ where $\mathbf{P}(h) \propto e^{-\mathbf{r}(h)}$. The final form of the loss function is

$$L(h) = \sum_{n \in N} l_h(x_n, y_n) + \mathbf{r}(h).$$

4 The Sampling Problem

Consider the case where the training set D was derived by downsampling a larger dataset D^+ . Formally, we say that D was generated from D^+ with a *sampling probability function* \mathbf{s} .

We limit ourselves to samplers $\mathbf{s} : X \times Y \rightarrow [0, 1]$ that consider each datapoint $(x, y) \in D^+$ independently.

This type of sampling is optimal for parallel computation because the sampler has no state other than it's inputs (x, y) . It does exclude some common sampling types, some of which are covered in Section 7.

When the sampling function is known, a posterior distribution can still be computed from D to learn which predictors are more likely to fit D^+ .

5 The General Solution

Start with the unnormalized Bayes rule in equation (4), but now consider that the posterior distribution and likelihood both depend on the sampling function \mathbf{s} .

$$\mathbf{P}(h|D, \mathbf{s}) \propto \mathbf{P}(D|h, \mathbf{s})\mathbf{P}(h)$$

Let $\mathbf{P}(x_n \in D^+)$ represent the probability that any given input x_n will appear in the unbiased dataset D^+ . This allows us to break down the likelihood as follows.

¹For a great treatment on loss functions and their various tradeoffs, see A Tutorial on Energy Based Learning by Lecun[7]. The work also discusses strategies for dealing with predictors where normalization over Y is not practical thus avoiding \hat{f}_h .

$$\mathbf{P}(D|h, s) = \prod_{n \in N} \mathbf{P}(x_n, y_n|h, s) = \prod_{n \in N} \mathbf{P}(x_n \in D^+) \mathbf{P}(y_n|x_n, h, s) \propto \prod_{n \in N} \mathbf{P}(y_n|x_n, h, s)$$

The term $\mathbf{P}(x_n \in D^+)$ is constant with respect to h and can therefore be dropped in the proportionality statement. We are left with calculating the expression $\mathbf{P}(y_n|x_n, h, s)$ which conditions on x_n , h and s . The following *generative description* is a useful tool in understanding how the sampling function factors into the generation of y_n .

1. Consider as given an input x_n , a predictor h , and a sampling function s . The predictor h encodes a probability distribution over Y through the function $f_h(x_n, y_n)$.
2. Sample from that probability distribution and make this a candidate for y_n , called y_n^*
3. Compute the sampling rate $s(x_n, y_n^*)$ and use that rate to probabilistically determine whether y_n^* is accepted.
 - (a) If it is accepted, return $y_n = y_n^*$.
 - (b) If it is not accepted, return to step 2 to generated another candidate.

The sampling function must eventually halt which means that it cannot be 0 for every potential candidate of any given x_n . This is a fair assumption because otherwise x_n would never appear in D to begin with.

We now use the generative description to produce a recursive equation for $\mathbf{P}(y_n|x_n, h, s)$.

We want to know the probability of selecting y_n . Let $y \in Y$ be the first candidate. If y is accepted, this probability is equal to 1 if $y_n = y$ and 0 otherwise, given by the indicator function $[y_n = y]$.

If y is not accepted, then we the probability reverts to the original value of $\mathbf{P}(y_n|x_n, h, s)$. Therefore, the probabily of ultimately accepting y_n comes to

$$P(y_n|cand = y, x_n, h, s,) = s(x_n, y) [y_n = y] + (1 - s(x_n, y))P(y_n|x_n, h, s).$$

If Y is discrete, the probability of selecting y as candidate in the first place is $\hat{f}_h(x_n, y)$. Use this to sum over the probabilities of selecting each possible candidate and setup a recursive equation.

$$\mathbf{P}(y_n|x_n, h, s) = \sum_{y \in Y} \hat{f}_h(x_n, y) (s(x_n, y) [y_n = y] + (1 - s(x_n, y))P(y_n|x_n, h, s)) \quad (6)$$

With algebraic manipulation documented in appendix A, we solve for $\mathbf{P}(y_n|x_n, h, s)$, reduce \hat{f} to f , and derive the formulas for both discrete and continuous² Y .

$$\mathbf{P}(y_n|x_n, h, s) = \frac{f_h(x_n, y_n)s(s_n, y_n)}{\sum_{y \in Y} f_h(x_n, y)s(x_n, y)} \quad \mathbf{P}(y_n|x_n, h, s) = \frac{f_h(x_n, y_n)s(s_n, y_n)}{\int_{y \in Y} f_h(x_n, y)s(x_n, y)} \quad (7)$$

²The continuous version of this argument requires more mathematical background but is completely analogous. It requires integrating over all candidates instead of taking the sum. Some care must be taken with the term for the indicator function for a rigorous argument, but ultimately it can be reduced in the same way.

The feasibility of computing the sum or integral term depends on the structure of Y , but is at least easy when Y is finite and small enough to be enumerated by machine.

5.1 Formula for Negative Log-Likelihood

Equation (7) provides all the tools needed to assign relative posterior probabilities to predictors. We derive a negative log likelihood loss function for selection or sampling starting with the relative Bayes formula.

$$\mathbf{P}(h|D, \mathbf{s}) \propto \prod_{n \in N} [\mathbf{P}(y_n|x_n, h, \mathbf{s})] \mathbf{P}(h)$$

Use equation (7) to get this in terms of f_h :

$$\mathbf{P}(h|D, \mathbf{s}) \propto \prod_{n \in N} \left[\frac{f_h(x_n, y_n) \mathbf{s}(x_n, y_n)}{\sum_{y \in Y} f_h(x_n, y) \mathbf{s}(x_n, y)} \right] \mathbf{P}(h)$$

Finally, derive a negative log likelihood loss function on h :

$$L(h) = \sum_{n \in N} \left[l_h(x_n, y_n) - \ln \mathbf{s}(x_n, y_n) + \ln \sum_{y \in Y} f_h(x_n, y) \mathbf{s}(x_n, y) \right] + \mathbf{r}(h)$$

6 Solution for Binary Logistic Regression

Binary logistic regression is a special case of the supervised learning problem in section 3 with the following additional properties:

1. It is a *binary classification* in that $Y = \{0, 1\}$.
2. The input space X is a list of real valued features. Let F denote a finite set of features, $X = \mathbb{R}^{|F|}$.
3. Each predictor $h \in H$ is parameterized by $h = (c, \mathbf{w})$ where $c \in \mathbb{R}$ is the intercept and $\mathbf{w} \in \mathbb{R}^{|F|}$ is an $|F|$ -dimensional vector of weights corresponding to each input feature.

Each hypothesis $(c, \mathbf{w}) \in H$ corresponds to the following probability distribution function:

$$f_{c, \mathbf{w}}(x_n, y_n) = \frac{e^{y_n \cdot (c + \mathbf{w} \cdot x_n)}}{1 + e^{c + \mathbf{w} \cdot x_n}}$$

Use equation (7) to get

$$\mathbf{P}(y_n|x_n, h, \mathbf{s}) = \frac{f_h(x_n, y_n)\mathbf{s}(s_n, y_n)}{\sum_{y \in Y} f_h(x_n, y)\mathbf{s}(x_n, y)} = \frac{e^{y_n \cdot (c + \mathbf{w} \cdot x_n)} \mathbf{s}(s_n, y_n)}{\mathbf{s}(x_n, 0) + e^{c + \mathbf{w} \cdot x_n} \mathbf{s}(x_n, 1)}.$$

These problems are often framed to focus on the probability of the *target condition* $y_n = 1$. This target condition is usually the rare event that triggered the decision to use biased sampling, and would correspond to the lion image in section 2.1 and the visit in section 2.2. We can solve for it as follows:

$$\mathbf{P}(y_n = 1|x_n, h, \mathbf{s}) = \frac{e^{c + \mathbf{w} \cdot x_n} \mathbf{s}(s_n, 1)}{\mathbf{s}(x_n, 0) + e^{c + \mathbf{w} \cdot x_n} \mathbf{s}(x_n, 1)} = \frac{e^{c + \mathbf{w} \cdot x_n}}{\mathbf{s}_r(x_n) + e^{c + \mathbf{w} \cdot x_n}} \quad \text{where} \quad \mathbf{s}_r(x_n) = \frac{\mathbf{s}(x_n, 0)}{\mathbf{s}(x_n, 1)} \quad (8)$$

Here, $\mathbf{s}_r : X \rightarrow [0, \infty]$ is the true-to-false *sample ratio* for each instance n . This ratio is all that is needed to correct for sampling bias in binary classification. Note that for $\mathbf{s}_r(x_n) = 1$, we are left with the original binary logistic regression formula.

For the prior, we can use a Gaussian distribution (aka L2, or ridge regression) with weight λ so:

$$\mathbf{r}(c, \mathbf{w}) = \frac{1}{2} \lambda (\mathbf{w} \cdot \mathbf{w})$$

Put this together and drop some constant factors to derive a negative log-likelihood loss function.

$$L(h) = \sum_{n \in N} (\ln(\mathbf{s}_r(x_n) + e^{c + \mathbf{w} \cdot x_n}) - y_n \cdot (c + \mathbf{w} \cdot x_n)) + \frac{1}{2} \lambda (\mathbf{w} \cdot \mathbf{w})$$

The derivative on a single weight \mathbf{w}_f or intercept c can be computed into the following simple form in order to perform gradient descent.

$$\frac{\partial}{\partial \mathbf{w}_f} L(h) = \sum_{n \in N} x_{n,f} \left(\frac{e^{c + \mathbf{w} \cdot x_n}}{\mathbf{s}_r(x_n) + e^{c + \mathbf{w} \cdot x_n}} - y_n \right) + \lambda \cdot w_f \quad \frac{\partial}{\partial c} L(h) = \sum_{n \in N} \left(\frac{e^{c + \mathbf{w} \cdot x_n}}{\mathbf{s}_r(x_n) + e^{c + \mathbf{w} \cdot x_n}} - y_n \right) \quad (9)$$

7 Future Work

7.1 Simple Random Sampling

Simple Random Sampling (SRS) is where an exact number of datapoints are retained in the dataset. The point-wise sampling function s can never guarantee a specified number of instances will be selected. Simple Random Sampling can furthermore be stratified by label to account for rare events. Because the sampling process is no longer independant by instance, deriving the formula is a more complex task.

7.2 Oversampling

The sampling function \mathbf{s} assumes that each datapoint is either included in the dataset or excluded. This is known as *undersampling*. We could allow for *oversampling*, where some instances are included in D multiple times.

Now instead of $\mathbf{s} : X \times Y \rightarrow [0, 1]$, the sampling function s returns a probability distribution over all natural numbers. This could be any probability distribution, but in practice this is often the *poisson distribution*. The poisson distribution with parameter λ places probability of multiplicity k at $\frac{\lambda^k e^{-\lambda}}{k!}$.

Bootstrap sampling methods rely on oversampling.

7.3 The Value of an Instance

The ability to quantify and quickly estimate the expected value that an instance will provide if included in D would allow engineers to deploy samplers that eliminate unnecessary work.

Perhaps insights on sampling can be borrowed from the field of *active learning*, where learners can directly choose inputs to query. Angluin[1] provides foundational analysis on the efficiencies that can be achieved with active strategies. Many such strategies have been deployed for statistical models, and for example Cohn et al.[3] propose selecting input data to minimize learner variance.

Ultimately, we do not know in general how much a datapoint will change the posterior distribution before it is included, but estimation techniques could be developed. The engineer must also consider the specific goals and tradeoffs of the project to inform the sampling strategy.

7.4 Ensemble Models

If models produced from several sampling functions can be combined in an *ensemble model*, they will likely perform better than in the singular approach. Because they can be run in parallel, and each part runs quickly with a small sample size, this provides real engineering benefits.

Appendices

A Solving for the General Formula

Here we start with equation (6) and go through the series of steps necessary to derive its final form in equation (7). Equation (6) begins as

$$\mathbf{P}(y_n|x_n, h, \mathbf{s}) = \sum_{y \in Y} \hat{f}_h(x_n, y) (\mathbf{s}(x_n, y) [y_n = y] + (1 - \mathbf{s}(x_n, y)) P(y_n|x_n, h, \mathbf{s})).$$

Break out the summation to get

$$\mathbf{P}(y_n|x_n, h, \mathbf{s}) = \sum_{y \in Y} \hat{f}_h(x_n, y) \mathbf{s}(x_n, y) [y_n = y] + \sum_{y \in Y} \hat{f}_h(x_n, y) (1 - \mathbf{s}(x_n, y)) P(y_n|x_n, h, \mathbf{s}).$$

In the first sum, the only non-zero addend is $y = y_n$. Therefore, we can replace y with y_n and remove the summation and indicator function. In the second sum, factor out $P(y_n|x_n, h, \mathbf{s})$ which does not contain summation index y .

$$\mathbf{P}(y_n|x_n, h, \mathbf{s}) = \hat{f}_h(x_n, y_n) \mathbf{s}(x_n, y_n) + P(y_n|x_n, h, \mathbf{s}) \sum_{y \in Y} \hat{f}_h(x_n, y) (1 - \mathbf{s}(x_n, y))$$

Collect the term $\mathbf{P}(y_n|x_n, h, \mathbf{s})$ and distribute $\hat{f}_h(x_n, y)$ in the remaining summation.

$$\begin{aligned} \mathbf{P}(y_n|x_n, h, \mathbf{s}) \left[1 - \sum_{y \in Y} \hat{f}_h(x_n, y) (1 - \mathbf{s}(x_n, y)) \right] &= \hat{f}_h(x_n, y_n) \mathbf{s}(x_n, y_n) \\ \mathbf{P}(y_n|x_n, h, \mathbf{s}) \left[1 - \sum_{y \in Y} \hat{f}_h(x_n, y) + \sum_{y \in Y} \hat{f}_h(x_n, y) \mathbf{s}(x_n, y) \right] &= \hat{f}_h(x_n, y_n) \mathbf{s}(x_n, y_n) \end{aligned}$$

Because we are using the normalized probability function \hat{f} , and from equation (1) we know that for all possible inputs x , $\sum_{y \in Y} \hat{f}_h(x, y) = 1$, we can replace the summation over \hat{f} and cancel as follows:

$$\begin{aligned} \mathbf{P}(y_n|x_n, h, \mathbf{s}) \left[1 - 1 + \sum_{y \in Y} \hat{f}_h(x_n, y) \mathbf{s}(x_n, y) \right] &= \hat{f}_h(x_n, y_n) \mathbf{s}(x_n, y_n) \\ \mathbf{P}(y_n|x_n, h, \mathbf{s}) \left[\sum_{y \in Y} \hat{f}_h(x_n, y) \mathbf{s}(x_n, y) \right] &= \hat{f}_h(x_n, y_n) \mathbf{s}(x_n, y_n) \end{aligned}$$

Finally, we note that because of the equation (1), the \hat{f}_h terms are simply a constant factor of same terms with the unnormalized f_h . Because this normalization factor appears on both sides of the equation, \hat{f}_h can be reduced to f through cancellation. Thus with one extra step of division, the final form can be given as (7).

$$\mathbf{P}(y_n|x_n, h, \mathbf{s}) = \frac{f_h(x_n, y_n) \mathbf{s}(x_n, y_n)}{\sum_{y \in Y} f_h(x_n, y) \mathbf{s}(x_n, y)} \quad (7)$$

If we assume that there was no sampling, or the sampling was uniform where $\mathbf{s}(x, y) = p$, equation (7) should reduce to the original predictor probability function $\hat{f}_h(x_n, y_n)$ as a sanity check.

$$\mathbf{P}(y_n|x_n, h, \mathbf{s}) = \frac{f_h(x_n, y_n)\mathbf{s}(x_n, y_n)}{\sum_{y \in Y} f_h(x_n, y)\mathbf{s}(x_n, y)} = \frac{f_h(x_n, y_n)p}{\sum_{y \in Y} f_h(x_n, y)p} = \hat{f}_h(x_n, y_n)$$

References

- [1] Angluin, D. (1988). Queries and concept learning. *Machine learning*, 2(4), 319-342.
- [2] Blais, B. S. (2014). *Statistical Inference for Everyone* (sie).
- [3] Cohn, D. A., Ghahramani, Z., & Jordan, M. I. (1996). Active learning with statistical models. *Journal of artificial intelligence research*, 4, 129-145.
- [4] Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, B. D. (2013). *Bayesian Data Analysis*. 3rd edition.
- [5] Hoffman, M. D., & Gelman, A. (2014). The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *J. Mach. Learn. Res.*, 15(1), 1593-1623.
- [6] King, G., & Zeng, L. (2001). Logistic regression in rare events data. *Political analysis*, 9(2), 137-163.
- [7] LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M., & Huang, F. (2006). A tutorial on energy-based learning. *Predicting structured data*, 1(0).
- [8] Maalouf, M., & Siddiqi, M. (2014). Weighted logistic regression for large-scale imbalanced and rare events data. *Knowledge-Based Systems*, 59, 142-148.
- [9] Maalouf, M., & Trafalis, T. B. (2011). Rare events and imbalanced datasets: an overview. *International Journal of Data Mining, Modelling and Management*, 3(4), 375-388.
- [10] Martin, O. (2016). *Bayesian analysis with python*. Packt Publishing Ltd.
- [11] Salvatier, J., Wiecki, T. V., & Fonnesbeck, C. (2016). Probabilistic programming in Python using PyMC3. *PeerJ Computer Science*, 2, e55.
- [12] Sklar, M. (2014). Fast MLE computation for the Dirichlet multinomial. *arXiv preprint arXiv:1405.0099*.
- [13] Sklar, M., Stewart, R., Li, R., Bakula, A., & Spears, E. (2020). U.S. Patent Application No. 16/405,481. [Section 0037]