

Frequentist, Conditionalist, and Their Relationship with Machine Learning

Lili Mou
moull12@sei.pku.edu.cn

Outline

Introduction

Pathologies of Frequentist

Bayesian Perspective

Machine Learning

Outline

Introduction

Pathologies of Frequentist

Bayesian Perspective

Machine Learning

Probability

Kolmogorov (1933):

- ▶ Nonnegative

$$p(x) \geq 0, \forall x$$

- ▶ Normalized

$$\sum_x p(x) = 1$$

- ▶ Finitely/countably additive

Let A_1, A_2, \dots, A_n be disjoint events,

$$p\left(\bigcup A_i\right) = \sum_i p(A_i)$$

Interpretations

- ▶ The probability that a toss of a coin gives the head
- ▶ The probability that it will rain tomorrow
- ▶ The probability that the speed of light lies in $2.9\text{--}3.1 \times 10^8 \text{ m/s}$
- ▶ **Frequentist** The limit of frequency provided that the number of samples goes to infinity (Recall the Law of Large Numbers)
- ▶ **Bayesian** The degree of ones subjective belief
 - ▶ Is belief necessarily a kind of probability?
 - ▶ Is belief admissible in scientific research? or even unavoidable?

Interpretations

- ▶ The probability that a toss of a coin gives the head
- ▶ The probability that it will rain tomorrow
- ▶ The probability that the speed of light lies in $2.9\text{--}3.1 \times 10^8 \text{ m/s}$
- ▶ **Frequentist** The limit of frequency provided that the number of samples goes to infinity (Recall the Law of Large Numbers)
- ▶ **Bayesian** The degree of ones subjective belief
 - ▶ Is belief necessarily a kind of probability?
 - ▶ Is belief admissible in scientific research? or even unavoidable?

“Some even argue that the frequency concept never applies, it being impossible to have an infinite sequence of i.i.d repetitions of any situation, except in a certain imaginary (subjective) sense.” [1]

Interpretations

- ▶ The probability that a toss of a coin gives the head
- ▶ The probability that it will rain tomorrow
- ▶ The probability that the speed of light lies in $2.9\text{--}3.1 \times 10^8 \text{ m/s}$
- ▶ **Frequentist** The limit of frequency provided that the number of samples goes to infinity (Recall the Law of Large Numbers)
- ▶ **Bayesian** The degree of ones subjective belief
 - ▶ Is belief necessarily a kind of probability?
 - ▶ Is belief admissible in scientific research? or even unavoidable?

“Some even argue that the frequency concept never applies, it being impossible to have an infinite sequence of i.i.d repetitions of any situation, except in a certain imaginary (subjective) sense.” [1]

“The subjectivist states his judgements, whereas the objectivist sweeps them under the carpet by calling assumptions knowledge, and he basks under the glorious objectivity of science.” (Good, 1973; see also [1].)

Bifurcation of the Two Schools

We have the data \mathcal{D} , and the model, parametrized by Θ .

$$\mathcal{D} \sim p_{\Theta}(\cdot)$$

What do we take expectation on (for learning, inference, etc)?

Frequentist: \mathcal{D} , because there is nothing random about Θ
No random, no cry

Bayesian: Θ , because \mathcal{D} is known
Everything unknown is a random variable.

Bifurcation of the Two Schools

We have the data \mathcal{D} , and the model, parametrized by Θ .

$$\mathcal{D} \sim p_{\Theta}(\cdot)$$

What do we take expectation on (for learning, inference, etc)?

Frequentist: \mathcal{D} , because there is nothing random about Θ
No random, no cry

Bayesian: Θ , because \mathcal{D} is known
Everything unknown is a random variable.

Who is right? Sufficiency + Weak condition \Rightarrow Bayesian analysis

Outline

Introduction

Pathologies of Frequentist

Bayesian Perspective

Machine Learning

Pathologies

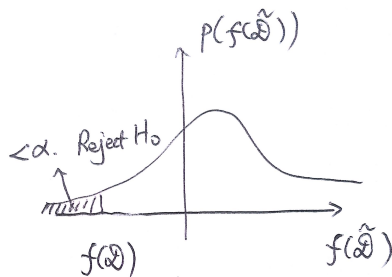
“Frequentist statistics exhibits various forms of weird and undesirable behaviors, known as pathologies.” [2]

- ▶ Hypothesis test
- ▶ Confidence interval

Hypothesis Test and p -Value

- ▶ $H_0 \leftrightarrow H_1$
- ▶ Data \mathcal{D}
- ▶ Test statistic: $f(\mathcal{D})$
- ▶ $p\text{-value}(\mathcal{D}) = \Pr\{f(\tilde{\mathcal{D}}) > f(\mathcal{D}) | \tilde{\mathcal{D}} \sim H_0\}$
- ▶ Reject H_0 , if $p\text{-value} < \alpha$
- ▶ Cannot reject H_0 , if $p\text{-value} \geq \alpha$

Tail Area Probability



First peculiar property: Because p -value relies on the tail area probability, “a hypothesis which may be true may be rejected because it has not predicted observable results which have not occurred.” (Jeffreys, 1961; see also [1].)

Who cares about nonoccurrence? But the disaster just begins. . .

Sample size

Assume $\mathcal{N}(\mu, 1)$

- ▶ $H_0 : \mu = 0 \leftrightarrow H_1 : \mu \neq 0$
- ▶ True probability: $\mathcal{N}(0.2, 1)$
- ▶ $\mathcal{D}_1 = \{0.2\}$, $\mathcal{D}_2 = 10000$ samples with mean 0.2

Given \mathcal{D}_1 , we cannot reject H_0 .

Given \mathcal{D}_2 , we do reject H_0 .

Sample size

Assume $\mathcal{N}(\mu, 1)$

- ▶ $H_0 : \mu = 0 \leftrightarrow H_1 : \mu \neq 0$
- ▶ True probability: $\mathcal{N}(0.2, 1)$
- ▶ $\mathcal{D}_1 = \{0.2\}$, $\mathcal{D}_2 = 10000$ samples with mean 0.2

Given \mathcal{D}_1 , we cannot reject H_0 .

Given \mathcal{D}_2 , we do reject H_0 .

For most scientific problems, the only thing that matters is the sample size, which is under control of researchers.

“In a recent survey, 58% of researchers admitted to having collected more data after looking to see whether the results were significant and 22% admitted to stopping an experiment early because they had found the result that they were looking for.” (Sanborn et al., 2014)

Interpretation

- ▶ Expected type I (false positive) error rate is at most α
- ▶ Nothing is said how often you err when you accept, or reject.

Interpretation

- ▶ Expected type I (false positive) error rate is at most α
- ▶ Nothing is said how often you err when you accept, or reject.

“This is sometimes interpreted as saying that frequentist hypothesis testing is very conservative, since it is unlikely to accidentally reject the null hypothesis. But in fact the opposite is the case: because this method only worries about trying to reject the null, it can never gather evidence in favor of the null, no matter how large the sample size. Because of this, p -values tend to overstate the evidence against the null, and are thus very ‘trigger happy.’ ” [2]

Interpretation

- ▶ Expected type I (false positive) error rate is at most α
- ▶ Nothing is said how often you err when you accept, or reject.

“This is sometimes interpreted as saying that frequentist hypothesis testing is very conservative, since it is unlikely to accidentally reject the null hypothesis. But in fact the opposite is the case: because this method only worries about trying to reject the null, it can never gather evidence in favor of the null, no matter how large the sample size. Because of this, p -values tend to overstate the evidence against the null, and are thus very ‘trigger happy.’ ” [2]

“It seems somewhat nonsensical, however, that we first deliberately formulate the problem wrong, and then in an *ad hoc* fashion explain the final results in more reasonable terms.” [1]

Confidence Interval

$$C_\alpha(\theta) = (l, u) : \Pr\{l(\tilde{\mathcal{D}}) \leq \theta \leq u(\tilde{\mathcal{D}}) | \tilde{\mathcal{D}} \sim \theta\} = 1 - \alpha$$

Counter-intuitive explanation:

- ▶ The confidence level (e.g., 95%) is **NOT** the probability that θ lies in the interval, given \mathcal{D} .
- ▶ It is the probability that the interval covers θ if we repeatedly draw datasets $\tilde{\mathcal{D}}$ (in addition to \mathcal{D} *per se*).
- ▶ However, we notice that \mathcal{D} is **KNOWN**. Who cares about nonoccurrence?

Twists and Turns

Pratt (1962):

- ▶ An engineer draws a random sample of electron tubes and measures the plate voltages under certain conditions with a very accurate voltmeter, accurate enough so that measurement error is negligible compared with the variability of the tubes.
- ▶ A statistician examines the measurements, which look normally distributed and vary from 75 to 99 volts with a mean of 87 and a standard deviation of 4. He makes the ordinary normal analysis, giving a confidence interval for the true mean.

Twists and Turns (2)

- ▶ Later he visits the engineer's laboratory, and notices that the voltmeter used reads only as far as 100, so the population appears to be “censored.” This necessitates a new analysis, if the statistician is orthodox.
- ▶ However, the engineer says he has another meter, equally accurate and reading to 1000 volts, which he would have used if any voltage had been over 100.
- ▶ This is a relief to the orthodox statistician, because it means the population was effectively uncensored after all.

Twists and Turns (3)

- ▶ But the next day the engineer telephones and says, “I just discovered my high-range voltmeter was not working the day I did the experiment you analyzed for me.”
- ▶ The statistician ascertains that the engineer would not have held up the experiment until the meter was fixed, and informs him that a new analysis will be required.
- ▶ The engineer is astounded. He says, “But the experiment turned out just the same as if the high-range meter had been working. I obtained the precise voltages of my sample anyway, so I learned exactly what I would have learned if the high-range meter had been available. Next you’ll be asking about my oscilloscope.”

Outline

Introduction

Pathologies of Frequentist

Bayesian Perspective

Machine Learning

Weak Conditionality

“Suppose a substance to be analyzed can be sent either to a laboratory in New York or a laboratory in California. The two labs seem equally good, so a fair coin is flipped to choose between them, which “heads” denoting that the lab in New York will be chosen. The coin is flipped and comes up tails, so the California lab is used. After a while, the experimental results come back and a conclusion and report must be developed. Should this conclusion take into account the fact that the coin could have been heads, and hence that the experiment in New York might have been performed instead?”

Common sense (and the conditional viewpoint) cries no, that only the experiment actually performed is relevant, but frequentist reasoning would call for averaging over all possible data, even the possible New York data.

Sufficiency + Weak Conditionality Principle \Rightarrow Bayesian Analysis

Outline

Introduction

Pathologies of Frequentist

Bayesian Perspective

Machine Learning

Linear Classification [3]

Assume

$$p(y = i | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$$

Given a set of training samples \mathbf{x}_i, y_i , we would like to predict $y_* = ?$ for a new sample \mathbf{x}_* .

Orthodox Frequentists' Viewpoint

Minimize the expected risk (loss) over \mathbf{x}_*



Minimize the empirical risk over \mathbf{x}_i, y_i (training samples)



Maximize the likelihood of \mathbf{x}_i, y_i

Maximum likelihood estimation

Training:

$$\mathbf{w}^* \leftarrow \arg \max_{\mathbf{w}} p(y_i | \mathbf{w}, \mathbf{x}_i)$$

Predicting:

$$y_* \leftarrow \arg \max_y p(y_* | \mathbf{x}_*, \mathbf{w}^*)$$

Deceitful Frequentists' Viewpoint

We somehow know that God created the world in a neat manner, and thus w is small.

We assume that $w \sim \mathcal{N}(0, \sigma)$ *a priori*.

We therefore maximize the posterior rather than likelihood.

Training:

$$\begin{aligned}w^* &\leftarrow \arg \max_w p(w|y_i, \mathbf{x}_i) \\ &= \arg \max_w p(w)p(y_i|w, \mathbf{x}_i)\end{aligned}$$

Predicting:

$$y_* \leftarrow \arg \max_y p(y_*|\mathbf{x}_*; w^*)$$

Gaussian prior $\Rightarrow \ell_2$ penalty

Laplacian prior $\Rightarrow \ell_1$ penalty

Bayesian Learning

- ▶ There does not exist \mathbf{w}^* .
- ▶ Rather, \mathbf{w} is a random variable that we have to marginalize out.
- ▶ Predictive density

$$p(y_*|\mathbf{y}) = \int d\mathbf{w} p(\mathbf{w}|\mathbf{y})p(y_*|\mathbf{w})$$

Bayesian Learning

- ▶ There does not exist \mathbf{w}^* .
- ▶ Rather, \mathbf{w} is a random variable that we have to marginalize out.
- ▶ Predictive density

$$p(y_*|\mathbf{y}) = \int d\mathbf{w} p(\mathbf{w}|\mathbf{y})p(y_*|\mathbf{w})$$

Warning: Formula menagerie ahead!

Bayesian Logistic Regression

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Bayesian Logistic Regression

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Prior $p(\mathbf{w})$, which is a \$64,000,000 question

Bayesian Logistic Regression

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Prior $p(\mathbf{w})$, which is a \$64,000,000 question

Likelihood

$$p(\mathbf{y}|\mathbf{w}) = \prod_{i=1}^m p(y^{(i)}|\mathbf{w})$$

Bayesian Logistic Regression

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Prior $p(\mathbf{w})$, which is a \$64,000,000 question

Likelihood

$$p(\mathbf{y}|\mathbf{w}) = \prod_{i=1}^m p(y^{(i)}|\mathbf{w}) = \prod_{i=1}^m \sigma(\mathbf{w}^T \phi^{(i)})^{t^{(i)}} (1 - \sigma(\mathbf{w}^T \phi^{(i)}))^{1-t^{(i)}}$$

Bayesian Logistic Regression

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Prior $p(\mathbf{w})$, which is a \$64,000,000 question

Likelihood

$$p(\mathbf{y}|\mathbf{w}) = \prod_{i=1}^m p(y^{(i)}|\mathbf{w}) = \prod_{i=1}^m \sigma(\mathbf{w}^T \phi^{(i)})^{t^{(i)}} (1 - \sigma(\mathbf{w}^T \phi^{(i)}))^{1-t^{(i)}}$$

Posterior

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{w})p(\mathbf{y}|\mathbf{w})}{p(\mathbf{y})}$$

Bayesian Logistic Regression

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Prior $p(\mathbf{w})$, which is a \$64,000,000 question

Likelihood

$$p(\mathbf{y}|\mathbf{w}) = \prod_{i=1}^m p(y^{(i)}|\mathbf{w}) = \prod_{i=1}^m \sigma(\mathbf{w}^T \phi^{(i)})^{t^{(i)}} (1 - \sigma(\mathbf{w}^T \phi^{(i)}))^{1-t^{(i)}}$$

Posterior

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{w})p(\mathbf{y}|\mathbf{w})}{p(\mathbf{y})} \propto_{\mathbf{w}} p(\mathbf{w})p(\mathbf{y}|\mathbf{w})$$

Bayesian Logistic Regression

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Prior $p(\mathbf{w})$, which is a \$64,000,000 question

Likelihood

$$p(\mathbf{y}|\mathbf{w}) = \prod_{i=1}^m p(y^{(i)}|\mathbf{w}) = \prod_{i=1}^m \sigma(\mathbf{w}^T \phi^{(i)})^{t^{(i)}} (1 - \sigma(\mathbf{w}^T \phi^{(i)}))^{1-t^{(i)}}$$

Posterior

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{w})p(\mathbf{y}|\mathbf{w})}{p(\mathbf{y})} \propto_{\mathbf{w}} p(\mathbf{w})p(\mathbf{y}|\mathbf{w}) = p(\mathbf{w}) \prod_{i=1}^m \sigma(\cdot)^{t^{(i)}} (1 - \sigma(\cdot))^{1-t^{(i)}}$$

Bayesian Logistic Regression

Let $\mathbf{y} \in \{0, 1\}^m$ denote the labels of training data ϕ_1, \dots, ϕ_m

Prior $p(\mathbf{w})$, which is a \$64,000,000 question

Likelihood

$$p(\mathbf{y}|\mathbf{w}) = \prod_{i=1}^m p(y^{(i)}|\mathbf{w}) = \prod_{i=1}^m \sigma(\mathbf{w}^T \phi^{(i)})^{t^{(i)}} (1 - \sigma(\mathbf{w}^T \phi^{(i)}))^{1-t^{(i)}}$$

Posterior

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{w})p(\mathbf{y}|\mathbf{w})}{p(\mathbf{y})} \propto_{\mathbf{w}} p(\mathbf{w})p(\mathbf{y}|\mathbf{w}) = p(\mathbf{w}) \prod_{i=1}^m \sigma(\cdot)^{t^{(i)}} (1 - \sigma(\cdot))^{1-t^{(i)}}$$

Predictive density

$$p(y_*|\mathbf{y}) = \int d\mathbf{w} p(y_*|\mathbf{w}) \cdot p(\mathbf{w}|\mathbf{y})$$
$$\propto_{y_*} \int d\mathbf{w} \sigma(\mathbf{w}^T \phi_*) \cdot p(\mathbf{w}) \prod_{i=1}^m \sigma(\cdot)^{t^{(i)}} (1 - \sigma(\cdot))^{1-t^{(i)}}$$

Intractability

- ▶ Posterior is intractable due to the normalizing factor.
- ▶ Predictive density is intractable due to the integral.

Intractability

- ▶ Posterior is intractable due to the normalizing factor.
- ▶ Predictive density is intractable due to the integral.

We have to resort to approximations

- ▶ Sampling methods
 - Stochastic, usually asymptotically correct, hard to scale
- ▶ Deterministic methods
 - “Do things wrongly and hope they work”

Laplace Approximation

- ▶ Fit a Gaussian at a mode
- ▶ The standard deviation is chosen such that ...
the second-order derivative of the log probability matches
- ☺ The first-order derivative is always 0 at a mode
- ☺ Scale free in representing the unnormalized measure
- ☹ Real variables only
- ☹ Only local properties captured, multi-mode distributions?

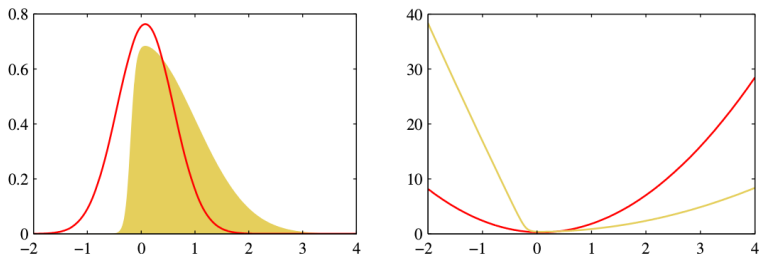


Figure 4.14 Illustration of the Laplace approximation applied to the distribution $p(z) \propto \exp(-z^2/2)\sigma(20z+4)$ where $\sigma(z)$ is the logistic sigmoid function defined by $\sigma(z) = (1 + e^{-z})^{-1}$. The left plot shows the normalized distribution $p(z)$ in yellow, together with the Laplace approximation curves centred on the mode z_0 of $p(z)$ in red. The right plot shows the negative logarithms of the corresponding curves.

Fitting a Gaussian

Let $p(z) = \frac{1}{Z} f(z)$ be a true distribution, where $Z = \int f(z) dz$

Step 1: Find a mode z_0 of $p(z)$, by gradient methods, say, satisfying

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0$$

Fitting a Gaussian

Let $p(z) = \frac{1}{Z} f(z)$ be a true distribution, where $Z = \int f(z) dz$

Step 1: Find a mode z_0 of $p(z)$, by gradient methods, say, satisfying

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0$$

Step 2: Consider a Taylor expansion of $\ln f(z)$ at z_0

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2} A (z - z_0)^2$$

where A is given by

$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

Taking the exponential,

$$f(z) \simeq f(z_0) \exp \left\{ - \frac{A}{2} (z - z_0)^2 \right\}$$

Fitting a Gaussian

Let $p(z) = \frac{1}{Z} f(z)$ be a true distribution, where $Z = \int f(z) dz$

Step 1: Find a mode z_0 of $p(z)$, by gradient methods, say, satisfying

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0$$

Step 2: Consider a Taylor expansion of $\ln f(z)$ at z_0

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2} A (z - z_0)^2$$

where A is given by

$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

Taking the exponential,

$$f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$

Step 3: Normalize to a Gaussian distribution

$$q(z) = \left(\frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$

Laplace Approximation for Multivariate Distributions

To approximate $p(\mathbf{z}) = \frac{1}{Z}f(\mathbf{z})$, where $\mathbf{z} \in \mathbb{R}^m$

We expand at mode \mathbf{z}_0

$$\ln f(\mathbf{z}) \simeq \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)$$

where $\mathbf{A} = -\nabla\nabla \ln f(\mathbf{z})\Big|_{\mathbf{z}=\mathbf{z}_0}$, Hessian of $\ln f(\mathbf{z})$, serving as the *precision matrix* in a Gaussian distribution

Taking the exponential, we obtain

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\}$$

Normalize it as a distribution, and then we have

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{m/2}} \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1})$$

Bayesian Logistic Regression: A Revisit in Earnest

Prior: Gaussian, which is natural¹

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

Posterior: $p(\mathbf{w} | \mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t} | \mathbf{w})$

$$\ln p(\mathbf{w} | \mathbf{t}) = -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) \quad \text{[prior]}$$

$$+ \sum_{i=1}^m \left\{ t^{(i)} \ln y^{(i)} + (1 - t^{(i)}) \ln(1 - y^{(i)}) \right\} \quad \text{[likelihood]}$$

+ const

$$\mathbf{S}_N = -\nabla \nabla \ln p(\mathbf{w} | \mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{i=1}^m y^{(i)}(1 - y^{(i)}) \phi_n \phi_n^T$$

Hence, the Laplace approximation to the posterior is

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{\text{MAP}}, \mathbf{S}_N)$$

¹Mathematicians always choose priors for the sake of convenience rather than approaching God.

Predictive Density

$$p(\mathcal{C}_1|\boldsymbol{\phi}_*, \mathbf{t}) = \int p(\mathcal{C}_1|\boldsymbol{\phi}_*, \mathbf{w})p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \simeq \int \sigma(\mathbf{w}^T \boldsymbol{\phi}_*)q(\mathbf{w}) d\mathbf{w}$$

$$p(\mathcal{C}_2|\boldsymbol{\phi}_*, \mathbf{t}) = 1 - p(\mathcal{C}_1|\boldsymbol{\phi}_*, \mathbf{t})$$

Plan

- ▶ Change it to a univariate integral
- ▶ Substitute sigmoid with a probit function, which is then convolved with a normal

$$\int \sigma(\cdot)\mathcal{N}(\cdot) d\cdot \simeq \int \Phi(\cdot)\mathcal{N}(\cdot) d\cdot = \Phi(\cdot) \simeq \sigma(\cdot)$$

The Dirac Delta Function

Let δ be the *Dirac delta* function, loosely thought of a function such that

- ▶ Gaussian distribution peaked at 0 with standard deviation $\rightarrow 0$

- ▶
$$\delta(x) = \begin{cases} +\infty, & \text{if } x = 0 \\ 0, & \text{otherwise} \end{cases}$$

δ function satisfies

$$\int \delta(a - x) f(a) da = f(x)$$

and specifically

$$\int \delta(a - x) da = 1$$

Deriving the Predictive Density

$$p(\mathcal{C}_1|\phi_*) \simeq \int \sigma(\mathbf{w}^T \phi_*) q(\mathbf{w}) d\mathbf{w} \quad \text{[Laplace approx.]}$$

$$\sigma(\mathbf{w}^T \phi) = \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da \quad \text{[Def. of } \delta \text{]}$$

$$\begin{aligned} p(\mathcal{C}_1|\phi_*) &= \int \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da q(\mathbf{w}) d\mathbf{w} \\ &= \int \int \sigma(a) \delta(a - \mathbf{w}^T \phi_*) q(\mathbf{w}) d\mathbf{w} da \\ &\stackrel{\Delta}{=} \int \sigma(a) p(a) da \end{aligned}$$

where

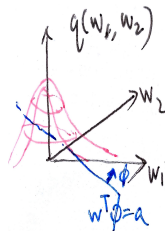
$$p(a) \stackrel{\Delta}{=} \int \delta(a - \mathbf{w}^T \phi_*) q(\mathbf{w}) d\mathbf{w}$$

We now argue that $p(a)$ is Gaussian

Deriving the Predictive Density (2)

$$\begin{aligned} p(a) &= \int \int \cdots \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) q(\mathbf{w}) dw_1 dw_2 \cdots dw_n \\ &= \int \int \cdots \int q(\tilde{\mathbf{w}}) dw_2 \cdots dw_n \end{aligned}$$

$\tilde{\mathbf{w}}$ is such that $\tilde{\mathbf{w}}^T \boldsymbol{\phi}_* = a$



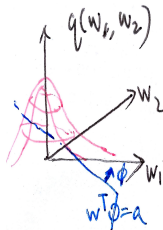
Deriving the Predictive Density (2)

$$\begin{aligned} p(a) &= \int \int \cdots \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) q(\mathbf{w}) d w_1 d w_2 \cdots d w_n \\ &= \int \int \cdots \int q(\tilde{\mathbf{w}}) d w_2 \cdots d w_n \end{aligned}$$

$\tilde{\mathbf{w}}$ is such that $\tilde{\mathbf{w}}^T \boldsymbol{\phi}_* = a$

We can also verify that

$$\begin{aligned} \int p(a) da &= \int \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) q(\mathbf{w}) d \mathbf{w} da \\ &= \int \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) q(\mathbf{w}) da d \mathbf{w} \\ &= \int q(\mathbf{w}) d \mathbf{w} \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}_*) da \\ &= 1 \end{aligned}$$



Deriving the Predictive Density (3)

$$\mu_a = \mathbb{E}[a] = \int p(a)a \, da = \int q(\mathbf{w})\mathbf{w}^T \boldsymbol{\phi}_* \, d\mathbf{w} = \mathbf{w}_{\text{MAP}}^T \boldsymbol{\phi}_*$$

$$\begin{aligned}\sigma_a^2 &= \text{var}[a] = \int p(a) \{a^2 - \mathbb{E}[a]^2\} \, da \\ &= \int q(\mathbf{w}) \{(\mathbf{w}^T \boldsymbol{\phi})^2 - (\mathbf{m}_N^T \boldsymbol{\phi})^2\} \, d\mathbf{w} \\ &= \boldsymbol{\phi}^T \mathbf{S}_N \boldsymbol{\phi}\end{aligned}$$

Thus

$$\begin{aligned}p(\mathcal{C}_1 | \mathbf{t}) &\simeq \int \sigma(a)p(a) \, da = \int \sigma(a)\mathcal{N}(a|\mu_a, \sigma_a^2) \, da \\ &\simeq \int \Phi(\lambda a)\mathcal{N}(a|\mu_a, \sigma_a^2) \, da = \Phi\left(\frac{\mu_a}{(\lambda^{-2} + \sigma_a^2)^{1/2}}\right) \simeq \sigma(\kappa(\sigma_a^2)\mu_a)\end{aligned}$$

$\lambda = \sqrt{\pi/8}$, $\kappa(\sigma_a^2) = (1 + \pi\sigma_a^2/8)^{-1/2}$, chosen such that the rescaled probit function has the same slope as sigmoid at the origin.

Hierarchical Bayes

What if we have parameters in the prior?

- ▶ Maximum likelihood estimation (Empirical Bayes, Type-II ML)
- ▶ Max *a posteriori*, assuming some prior on the hyper-parameters
- ▶ Full Bayesian treatment: Marginalize out all unknown variable!

Take-Home Messages

- ▶ Frequentist takes expectation on (known) data while conditioning on (unknown) θ
- ▶ Conditionalist (Bayesian) takes expectation on (unknown) θ while conditioning on (known) data
- ▶ Bayes Hierarchy
 - ▶ MLE
 - ▶ Max *a posteriori*
 - ▶ Empirical Bayes
 - ▶ Max *a posteriori* estimation of hyper-parameters
 - ▶ Full Bayes
- ▶ Bayesian treatment is fundamentally correct by computationally non-trivial.

References

- [1] James O. Berger, *Statistical Decision Theory and Bayesian Analysis (2 ed.)*, Springer-Verlag, 1985.
- [2] Kevin P. Murphy, *Machine Learning: A Probability Perspective*, MIT Press, 2012.
- [3] Christopher M. Bishop *Pattern Recognition and Machine Learning*, Springer, 2006.