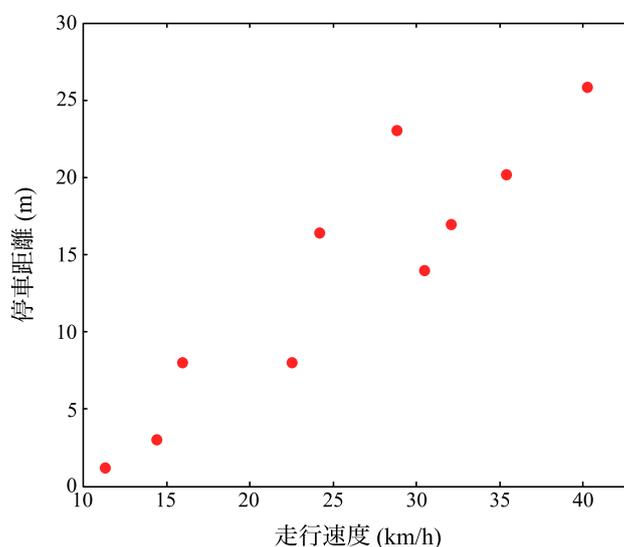


## 損失関数

ここでは、損失関数 (cost function) と呼ばれる関数について見ていきます。ここでの解説は、I. Goodfellow, Y. Bengio, and A. Courville, “Deep Learning” (2016) の 6.2.1 節の内容や数式を噛み砕いて、初めての人にも分かりやすくしたものです。関数の定義や記法 (notation) は、基本的にこの教科書に従っていますので、必要に応じて教科書の方も参照して下さい。

### ♣ 機械学習における損失関数とは？

具体的な数式の説明に入る前に、ざっくりとここで扱う損失関数とはどのようなものなのか、そのイメージを説明しておきます。前回の「条件付き対数尤度と平均二乗誤差」の解説において、車の走行速度と停車距離に関する次のデータを考えました。



もう1度、この10点のデータと同じ走行速度に対する停車距離を測定したとしましょう。上と完璧に同じデータが得られるでしょうか？ そのようなことは稀でしょう。実際の測定には誤差がつきものですから、データは揺らぎます。このような揺らぎを表現するための道具として確率分布が導入されます。統計の分野 (そして統計をベースにしている機械学習の分野) では、「データはある確率分布から生成される」と考えるのが基本です。つまり、上のデータも、走行速度を表す変数を  $x$ 、停車距離を表す変数を  $y$  とした時に、何らかの (条件付き) 確率分布  $p(y|x)$  が背後に存在し、10個のデータはこの確率分布に従って生成されたものと解釈されます。

今回考えるのは、与えられたデータ (上の例では10個の走行速度と停車距離のデータのセット) から、どのようにして背後にあるこの確率分布  $p(y|x)$  を得るか、ということです。もし仮に、完璧な  $p(y|x)$  を知ることができれば、任意の走行速度に対する停車距離を計算で簡単に予測することができるようになるため、非常に便利です。しかし、世の中そんな単純でなく、このよ

うなデータを生成する「真の確率分布」 $p(y|x)$ は、一般に非常に複雑で取り扱いが困難(不可能)であると考えられます。そこで、この真の確率分布 $p(y|x)$ をなるべく高い精度で再現でき、かつ比較的シンプルで取り扱いが容易な「確率モデル」 $p_{\text{model}}(y|x; \theta)$ を代わりに構築することを目指します。例えば、次のような Gauss 分布を採用するのは1つの手です<sup>1</sup>。

$$p_{\text{model}}(y|x; \theta) = \mathcal{N}(y; \theta_0 + \theta_1 x, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{[y - (\theta_0 + \theta_1 x)]^2}{2\sigma^2}\right)$$

$\theta$ はモデルを特徴づけるパラメータのベクトル(配列)で、この Gauss 分布のモデルの場合には、分布の平均を与える1次関数の切片 $\theta_0$ と傾き $\theta_1$ と、分布の分散 $\sigma^2$ の3つとなります。

$$\theta = (\theta_0, \theta_1, \sigma^2)$$

もちろん、実際のデータは、こんなに単純な確率分布に従って生成されてはいないでしょう。しかし、パラメータ $\theta$ を適切に調整して、確率モデル $p_{\text{model}}(y|x; \theta)$ を真の確率分布 $p(y|x)$ に十分近づけることが出来れば、高い精度での予測が可能となり実用上の問題はほとんどありません。機械学習における学習(訓練)とは、ざっくり言えば、このパラメータ $\theta$ を適切に調整して、確率モデル $p_{\text{model}}(y|x; \theta)$ を真の確率分布 $p(y|x)$ に近づける作業のことを指します。

このような学習を行うためには、確率モデルが真の確率分布からどのくらいズレているかを測る尺度が必要となります。この分布のズレ(=損失)を表す尺度が「**損失関数 (cost function)**」なのです。機械学習では、学習によりこの損失関数をなるべく小さく(最小化)することで、精度の高い予測ができるようになるのです。

### ♣ カルバック・ライブラー情報量と交差エントロピー

では、そのような分布のズレを表す損失関数とは、具体的にはどのような関数なのでしょうか。一般に、2つの確率分布 $P(x)$ と $Q(x)$ のズレ(離れ具合)を表す量としては、次の**カルバック・ライブラー情報量 (Kullback-Leibler divergence)**と呼ばれる量が知られています<sup>2</sup>。

$$D_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P}[\log P(x) - \log Q(x)]$$

上の式で出てくる $\mathbb{E}_{x \sim P}[f(x)]$ は、確率分布 $P(x)$ に関する $f(x)$ の期待値を表しています。期待値は「Bernoulli 分布」のところでも少し説明したのですが、この記号は初めて出てくるので、

<sup>1</sup>同じようなことを「条件付き対数尤度と平均二乗誤差」の解説でも行っているのですが、参考にして下さい。

<sup>2</sup>縦棒1本だと条件付き確率と紛らわしいので、 $P$ と $Q$ の間に2本縦棒を入れるのが慣例です。

もう1度おさらいしておきましょう。 $\mathbb{E}_{x \sim P}[f(x)]$  を具体的に数式で表すと、

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_i P(x_i) f(x_i)$$

となります<sup>3</sup>。確率変数  $x$  の取りうる全ての値  $x_i$  に対して、関数の値  $f(x_i)$  と  $x_i$  が現れる確率  $P(x_i)$  を掛け合わせたものを足す、という量になっています<sup>4</sup>。ちょっと分かりづらいかもしれないので具体例を考えましょう。次のようなゲームを考えます。100円支払うと、サイコロ1個を1回振ることができ、サイコロの目に応じて次の金額がもらえます。

サイコロの目 $x$	1	2	3	4	5	6
賞金 $f(x)$	20円	50円	100円	100円	150円	150円

出るサイコロの目を表す確率変数を  $x$  として、 $x_1$  を1の目、 $x_2$  を2の目、 $\dots$  と表すことにしましょう。賞金は出る目の数によって変化しますから、 $x$  の関数となっていて  $f(x)$  と表します。

$$f(x_1) = 20 \text{円}, \quad f(x_2) = 50 \text{円}, \quad \dots$$

そして出る目の確率は、(イカサマサイコロでなければ) どの目も  $\frac{1}{6}$  となるため、

$$P(x_1) = \frac{1}{6}, \quad P(x_2) = \frac{1}{6}, \quad \dots$$

となります。この時、貰える賞金の期待値  $\mathbb{E}_{x \sim P}[f(x)]$  を計算してみると、

$$\begin{aligned} \mathbb{E}_{x \sim P}[f(x)] &= \sum_{i=1}^6 P(x_i) f(x_i) \\ &= \frac{1}{6} \times 20 + \frac{1}{6} \times 50 + \frac{1}{6} \times 100 + \frac{1}{6} \times 100 + \frac{1}{6} \times 150 + \frac{1}{6} \times 150 = 95 \text{円} \end{aligned}$$

となります。つまり、このゲームで賞金として貰えると期待される金額は95円であり、参加費が100円であることを考えると、やや参加者が損するようなゲームになっていることが分かります。

ここまではサイコロの目のように、離散的(とびとび)な値を取る確率変数  $x$  を想定してきましたが、重さ、温度、速度などのように連続的な値を取る変数  $x$  を考える場合も当然あります。このような場合には、取りうる値が無限にあるため、先ほどのように総和記号を用いて期待値を表現することはできません。このような場合には、「積分」と呼ばれる計算によって期待値を

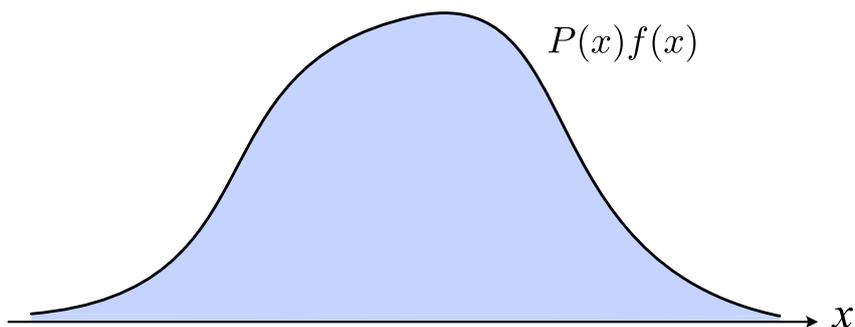
<sup>3</sup> $\Sigma$  は総和を表す記号です。softmax 関数のところで説明しています。

<sup>4</sup>一般に  $x$  の取りうる値の数がいくつあるか分からないので、和の範囲は書かず、単に  $\Sigma_i$  としています

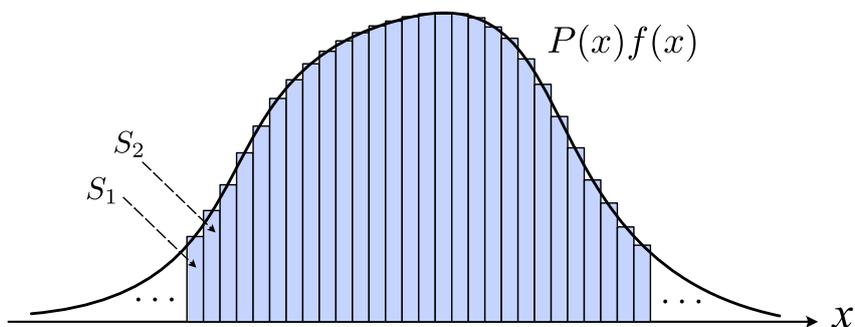
表すことができます。

$$\mathbb{E}_{x \sim P}[f(x)] = \int P(x)f(x) dx$$

積分は、総和  $\Sigma_i$  の計算を、 $x$  が連続的な値を取る場合に拡張した、いわば「特殊な総和」です。一説によると、この積分を表すニョロ  $\int$  は、総和を表す  $\Sigma$  の記号を引き伸ばして変形(進化?)させたものと言われています。具体的に積分とは「面積」を求める計算です。例えば、上の期待値の式での積分は、図形的には  $P(x)f(x)$  という関数と  $x$  軸とで挟まれた青い領域の面積を求めることに対応しています。



長方形や三角形の面積は、小学校で習う有名な公式を使えば簡単に求められますが、上の図のようなイビツな領域の面積を求めるためには、積分が必要となるのです。積分の基本的な考え方は、このイビツな図形を細長い長方形の短冊に分解して足し合わせるというものです。



短冊は長方形なので、その面積は縦 × 横で簡単に計算できます。ですから、イビツな図形でも細長い長方形の短冊に分けて、それぞれの短冊の面積  $S_i$  を計算して、各々の面積を全て足し合わせれば、面積を求めることができるのです<sup>5</sup>。その意味で、

$$\int P(x)f(x) dx = \sum_i S_i$$

になっており、確かに積分  $\int$  は、総和  $\Sigma$  を連続的な変数の場合に拡張したことになっているこ

<sup>5</sup>正確に面積を求めるためには、短冊は物凄く細かくする必要があります。

とが何となくわかるかと思います。具体的にこのような積分をどのように計算するかは、高校数学の範囲となります。ただし、特に具体的な積分の計算が分からなくても以下の内容は十分理解できるので、分からない方は、ひとまず、総和記号  $\Sigma$  を連続的な変数の場合に拡張したものの程度に思っておけば十分です。

期待値に関しては、いくつか重要な性質が成り立ちます。以下で使う性質として「期待値の和と差」に関する性質を確認しておきましょう。期待値に関しては、次の式が成り立ちます。

$$\mathbb{E}_{x \sim P}[f(x) + g(x)] = \mathbb{E}_{x \sim P}[f(x)] + \mathbb{E}_{x \sim P}[g(x)]$$

$$\mathbb{E}_{x \sim P}[f(x) - g(x)] = \mathbb{E}_{x \sim P}[f(x)] - \mathbb{E}_{x \sim P}[g(x)]$$

つまり、2つの関数の和(差)の期待値は、それぞれの関数の期待値の和(差)になっています。このことは、実際に総和(または積分)の形で期待値を表現してみると確かめられます。 $x$ が離散的な値を取る場合には、

$$\begin{aligned} \mathbb{E}_{x \sim P}[f(x) + g(x)] &= \sum_i P(x_i)[f(x_i) + g(x_i)] \\ &= \sum_i [P(x_i)f(x_i) + P(x_i)g(x_i)] \\ &= \sum_i P(x_i)f(x_i) + \sum_i P(x_i)g(x_i) = \mathbb{E}_{x \sim P}[f(x)] + \mathbb{E}_{x \sim P}[g(x)] \end{aligned}$$

といった感じで示せます。連続的な値を取る場合も同じ要領です。

だいぶ話が逸れてしまいましたが、期待値の確認をしたところで、再び Kullback-Leibler 情報量の式を眺めてみましょう<sup>6</sup>。

$$D_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P}[\log P(x) - \log Q(x)]$$

期待値の差に関する性質を用いれば、

$$D_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P}[\log P(x)] - \mathbb{E}_{x \sim P}[\log Q(x)]$$

ということになり、 $\log P(x)$  と  $\log Q(x)$  の確率分布  $P(x)$  に関する期待値の差を表していること

---

<sup>6</sup> $\log$  の底は何でもいいですが、普通はネイピア数  $e$  の自然対数を考えます。

が分かります。上の式の1項目と2項目の量にはそれぞれ名前がついています。1項目の

$$\mathbb{E}_{x \sim P}[\log P(x)] = \begin{cases} \sum_i P(x_i) \log P(x_i) & \text{(離散確率分布の場合)} \\ \int P(x) \log P(x) dx & \text{(連続確率分布の場合)} \end{cases}$$

は、「**エントロピー (entropy)**」と呼ばれる、確率分布  $P(x)$  の乱雑さ(予測のしにくさ)を表す量となっています。そして2項目の

$$\mathbb{E}_{x \sim P}[\log Q(x)] = \begin{cases} \sum_i P(x_i) \log Q(x_i) & \text{(離散確率分布の場合)} \\ \int P(x) \log Q(x) dx & \text{(連続確率分布の場合)} \end{cases}$$

は、「**交差エントロピー (cross entropy)**」と呼ばれます<sup>7</sup>。ですから、Kullback-Leibler 情報量  $D_{\text{KL}}(P||Q)$  は、分布  $P(x)$  のエントロピーから、分布  $Q(x)$  との交差エントロピーを引いた量となっています。

このセクションの最後に、具体的な2つの分布を考えて、確かに Kullback-Leibler 情報量が2つの分布の離れ具合を表す量になっていることを実感しておきましょう。ここでは例として、2つの Gauss 分布  $P(x)$ 、 $Q(x)$  の間の Kullback-Leibler 情報量を計算してみます<sup>8</sup>。

$$P(x) = \mathcal{N}(x; \mu_p, \sigma_p^2) = \frac{1}{\sqrt{2\pi}\sigma_p} \exp\left(-\frac{(x - \mu_p)^2}{2\sigma_p^2}\right)$$

$$Q(x) = \mathcal{N}(x; \mu_q, \sigma_q^2) = \frac{1}{\sqrt{2\pi}\sigma_q} \exp\left(-\frac{(x - \mu_q)^2}{2\sigma_q^2}\right)$$

この場合の Kullback-Leibler 情報量は、次のようになります。

$$\begin{aligned} D_{\text{KL}}(P||Q) &= \int P(x) \log P(x) dx - \int P(x) \log Q(x) dx = \int P(x) \log \frac{P(x)}{Q(x)} dx \\ &= \frac{1}{\sqrt{2\pi}\sigma_p} \int \exp\left(-\frac{(x - \mu_p)^2}{2\sigma_p^2}\right) \left[ \log\left(\frac{\sigma_q}{\sigma_p}\right) - \frac{(x - \mu_p)^2}{2\sigma_p^2} + \frac{(x - \mu_q)^2}{2\sigma_q^2} \right] dx \end{aligned}$$

かなり複雑な式変形に見えるかもしれませんが、単に、 $\log$  の性質を使っているだけです<sup>9</sup>。

<sup>7</sup>日本語でも「クロスエントロピー」と呼ぶこともありますが、ここでは交差エントロピーと呼びます。

<sup>8</sup>Gauss 分布については、以前の解説プリントを参照して下さい。

<sup>9</sup>-1 乗に関しては、softmax 関数の解説を見て下さい。

$$\log a - \log b = \log a + \log b^{-1} = \log a + \log \frac{1}{b} = \log \frac{a}{b}$$

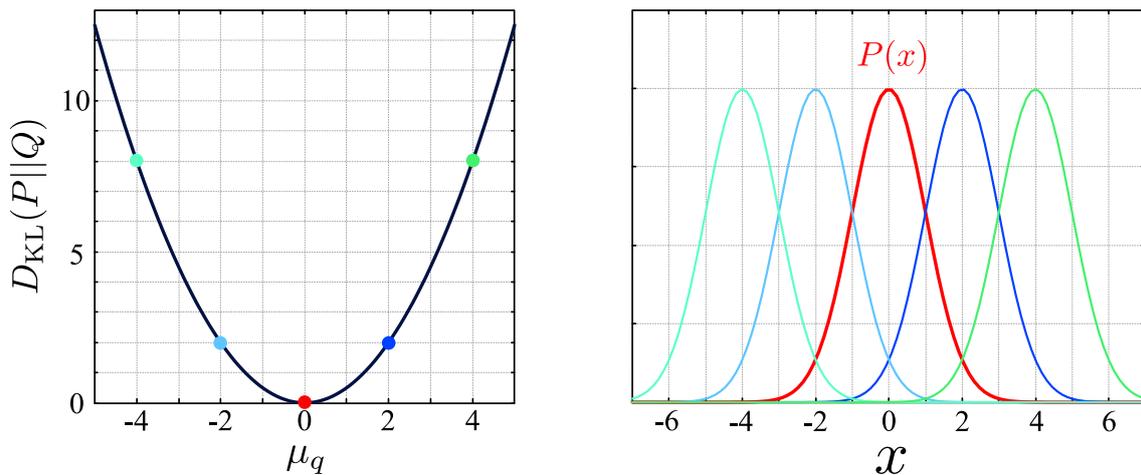
$$\log e^x = \log_e e^x = x \log_e e = x$$

ここでの対数logの底は、ネイピア数 $e$ であるとして計算しています。Gauss分布同士のKullback-Leibler情報量の計算に出てくる積分は簡単に計算することができて、計算結果は次のようになります。

$$D_{\text{KL}}(P||Q) = \log \frac{\sigma_q}{\sigma_p} + \frac{\sigma_p^2 + (\mu_p - \mu_q)^2}{2\sigma_q^2} - \frac{1}{2}$$

基本的な積分の計算は知っていて、さらにここでの積分の詳細が気になる方は「ガウス積分」などのキーワードで検索すれば色々出てくるので参考にして下さい。ただし単なる計算なので、ここでは結果を受け入れてもらえればそれで結構です<sup>10</sup>。

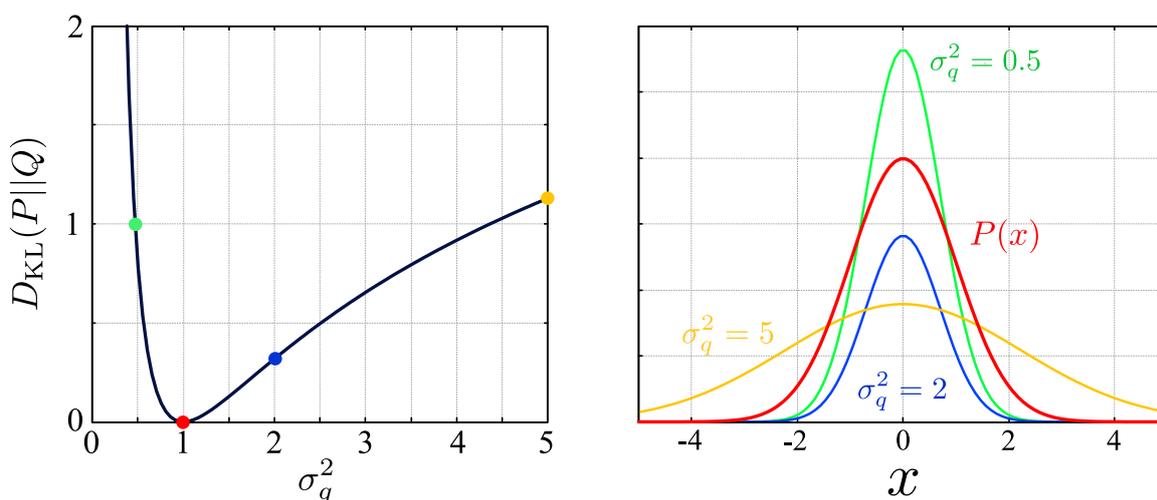
以下では、 $P(x)$ の分布の平均を $\mu_p = 0$ 、 $\sigma_p^2 = 1$ として、 $Q(x)$ の分布の平均 $\mu_q$ と分散 $\sigma_q^2$ を変化させた時の、Kullback-Leibler情報量の値の振る舞いを見ていきましょう。まず、平均 $\mu_q$ を変化させると次のようになります。



$\mu_q = 0$ の時、2つの分布は完璧に重なり、この時  $D_{\text{KL}}(P||Q) = 0$  となります。 $\mu_q$ を大きくしていくと、2つの分布は離れていき、それに伴って  $D_{\text{KL}}(P||Q)$ の値も大きくなっていく様子が見て取れます。

次に、分散 $\sigma_q^2$ を変化させた時の  $D_{\text{KL}}(P||Q)$ の変化を調べると、次のようになります。

<sup>10</sup>Wolfram Alpha (<https://www.wolframalpha.com>) という検索エンジンに式を入力すれば、積分の計算を一瞬でやってくれます。何か別の分布で計算してみたいけど、積分の計算が分からないという時に活用して下さい。



やはりこちらも、2つの分布が完全に一致する  $\sigma_q^2 = 1$  の時に  $D_{\text{KL}}(P||Q) = 0$  となり、そこから2つの分布がズレていくと  $D_{\text{KL}}(P||Q)$  の値も大きくなっていく様子が見て取れます。このことから、確かに Kullback-Leibler 情報量は、2つの分布の離れ具合を表す量になっていることが実感できるかと思います。

### ♣ 損失関数としての交差エントロピー (1)

さて、Kullback-Leibler 情報量が2つの分布の離れ具合を表す量であることが分かったため、この量は損失関数として用いることができます。入力値  $\mathbf{x}$  が与えられた時の出力値  $\mathbf{y}$  を生成する真の確率分布を  $p(\mathbf{y}|\mathbf{x})$  とします。そして、出力値  $\mathbf{y}$  を予想するためにパラメータ  $\boldsymbol{\theta}$  を適切に調整して構築していく確率モデルを  $p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  とします。

この2つの確率分布  $p(\mathbf{y}|\mathbf{x})$  と  $p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  を近づけることが、学習の目標です。条件付き確率であっても先ほどと同様に Kullback-Leibler 情報量を考えることができ、2つの分布の離れ具合は、

$$\begin{aligned} D_{\text{KL}}(\mathbf{p}||\mathbf{p}_{\text{model}}) &= \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [\log p(\mathbf{y}|\mathbf{x}) - \log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})] \\ &= \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [\log p(\mathbf{y}|\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [\log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})] \end{aligned}$$

と表すことができます。ここで1項目 (エントロピー) をみると、そこには  $p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  が一切含まれていないことが分かります。つまり、学習においてパラメータ  $\boldsymbol{\theta}$  を調整する際に、この部分は学習には何も影響しない、単なる定数になります。言い換えると、 $\boldsymbol{\theta}$  を変更、すなわちモデル  $p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  を変化させてもこの部分は全く変化しないため、学習の際には計算する必要がない部分となっています。計算する量はなるべく少ない方が効率的に学習できますから、

通常は2項目の

$$-\mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [\log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})] = \begin{cases} \sum_i p(\mathbf{y}_i|\mathbf{x}) \log p_{\text{model}}(\mathbf{y}_i|\mathbf{x}; \boldsymbol{\theta}) & \text{(離散確率分布の場合)} \\ \int p(\mathbf{y}|\mathbf{x}) \log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) d\mathbf{y} & \text{(連続確率分布の場合)} \end{cases}$$

の交差エントロピーの部分のみを、モデルと真の分布との差を表す損失関数として用います。ただし、交差エントロピーをそのまま損失関数に使おうとすると問題が生じます。上の量は、 $\mathbf{y}$  に関しては総和(または積分)が取られて、全ての場合が足し合わされているため、 $\mathbf{y}$  に依存することはありません。ただし、 $\boldsymbol{\theta}$  や  $\mathbf{x}$  を変化させると、その値は変化します。 $\boldsymbol{\theta}$  に依存することは全然問題でなく、むしろ  $\boldsymbol{\theta}$  を変化させてこの量を小さくすることで学習を行います。では、 $\mathbf{x}$  の値はどのようにすれば良いのでしょうか？ 入力値  $\mathbf{x}$  の値はサンプルによってバラバラで、そのどれを使えば良いかは明らかではありません。そこで、 $\mathbf{x}$  も確率分布に従って生成されるデータと考えると、それによって期待される交差エントロピーの値を考えることになります。この操作を理解するためには、結合確率、条件付き確率、周辺確率といった確率に関する基本事項を理解する必要がありますので、以下で見ていきましょう。

### 色々な確率 (結合確率、条件付き確率、周辺確率)

2つの確率変数  $x$ 、 $y$  が絡む確率には、様々な種類があります。例として、グー、チョキ、パーをランダムで出すジャンケンマシーンを2回動かした時に、グーがでる回数を  $x$ 、チョキがでる回数を  $y$  という確率変数で表すことにしましょう。 $x$  は0 (1回もグーが出ない) か1か2の値を取ります ( $y$  も同じです)。

この時、「グーが1回 ( $x = 1$ )、チョキが1回 ( $y = 1$ ) となるような確率」のような、 $x$  と  $y$  の値の組み合わせが生じる確率のことを「**結合確率 (joint probability)**」と呼び<sup>11</sup>、 $p(x, y)$  で表現します。例えば、2回ともグーが出て1回もチョキが出ないような確率は、

$$p(x = 2, y = 0) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

となります。一方、グーが2回、チョキが1回ということは、2回しかマシンは動かさないのでありえないため、

$$p(x = 2, y = 1) = 0$$

となります。これらを表にまとめると (統合確率分布といいます)、次のようになります。

---

<sup>11</sup>同時確率と呼ぶこともあります

$x \backslash y$	0	1	2
0	$\frac{1}{9}$	$\frac{2}{9}$	$\frac{1}{9}$
1	$\frac{2}{9}$	$\frac{2}{9}$	0
2	$\frac{1}{9}$	0	0

場合によっては、 $y$ の値は何でも良いけど  $x$ が0となる確率が知りたい、ということがあります。このような確率を「**周辺確率 (marginal probability)**」と呼び、 $p(x=0)$ で表現することになります。 $x=0$ の時、 $y$ は0~2の3通りの場合があるため、それらの確率を足し合わせると、

$$p(x=0) = p(x=0, y=0) + p(x=0, y=1) + p(x=0, y=2) = \frac{1}{9} + \frac{2}{9} + \frac{1}{9} = \frac{4}{9}$$

と計算できます。一般に、周辺分布は総和記号を使えば

$$p(x) = \sum_{i=1}^m p(x, y = y_i), \quad p(y) = \sum_{i=1}^m p(x = x_i, y)$$

と表せます。この周辺分布を先程の統合確率分布の表に書き足すと次のようになります。

$x \backslash y$	0	1	2	$p(x)$
0	$\frac{1}{9}$	$\frac{2}{9}$	$\frac{1}{9}$	$\frac{4}{9}$
1	$\frac{2}{9}$	$\frac{2}{9}$	0	$\frac{4}{9}$
2	$\frac{1}{9}$	0	0	$\frac{1}{9}$
$p(y)$	$\frac{4}{9}$	$\frac{4}{9}$	$\frac{1}{9}$	

右端にあるのが  $x$ に関する周辺確率  $p(x)$ 、下端にあるのが  $y$ に関する周辺確率  $p(y)$  となっています。お気づきだと思いますが、このように書いた時に表の周辺にくることから、「周辺」確率という名前がついています。

最後に、これは何度もこれまでに出てきましたが、確率変数  $x$  の値に依存して確率変数  $y$  の確率が決まる場合、「**条件付き確率 (conditional probability)**」と呼び、その分布を  $p(y|x)$  と書きます。この条件付き分布と統合確率、周辺確率の間には次の関係式が成り立ちます。

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

これは条件付き確率の意味を考えれば納得できます。  $p(y = b|x = a)$  は、  $x = a$  となる確率のうち、  $y = b$  である可能性が占める割合を意味します。すなわち、統合確率  $p(x = a, y = b)$  が周辺分布  $p(x = a)$  に対して占める割合によって求めることができるため、上の式のように表されます。

さて、確率の基本事項を確認したところで、元の話に戻りましょう。2つの確率分布  $p(\mathbf{y}|\mathbf{x})$  と  $p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  の交差エントロピー

$$-\mathbb{E}_{\mathbf{y} \sim p} [\log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})]$$

を、特定の  $\mathbf{x}$  への依存しないようにするために、確率分布  $p(\mathbf{x})$  についての期待値を考えましょう。ここで、  $p(\mathbf{x})$  は先ほど説明した  $\mathbf{x}$  の周辺分布を表します。離散確率分布の場合に、この期待値を計算してみると次のようになります。

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ -\mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [\log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})] \right] &= \sum_j \left[ p(\mathbf{x}_j) \left( -\sum_i p(\mathbf{y}_i|\mathbf{x}_j) \log p_{\text{model}}(\mathbf{y}_i|\mathbf{x}_j; \boldsymbol{\theta}) \right) \right] \\ &= \sum_i \sum_j \left[ -p(\mathbf{x}_j) p(\mathbf{y}_i|\mathbf{x}_j) \log p_{\text{model}}(\mathbf{y}_i|\mathbf{x}_j; \boldsymbol{\theta}) \right] \\ &= \sum_i \sum_j \left[ -p(\mathbf{x}_j, \mathbf{y}_i) \log p_{\text{model}}(\mathbf{y}_i|\mathbf{x}_j; \boldsymbol{\theta}) \right] \end{aligned}$$

難しい計算に感じるかもしれませんが、落ち着いて各ステップを考えればそこまで難しくありません。2つの期待値  $\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})}$  と  $\mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})}$  を表す総和記号の添字を  $i$  と  $j$  で区別していることに注意してください。プログラミングに慣れている方なら直ぐにわかると思いますが、これは2重ループになっているため、内側のループカウンタ変数と外側のループカウンタ変数が同じだとマズいですよね。そこで  $i$  と  $j$  という異なる添字 (カウンタ変数) を用いて区別しています。また、2行目から3行目に移るところで、種々の確率の間に成り立つ関係式

$$p(\mathbf{x}_j) p(\mathbf{y}_i|\mathbf{x}_j) = p(\mathbf{x}_j, \mathbf{y}_i)$$

を用いています。ややこしいのでもう1度確認しておく、  $p(\mathbf{x}_j)$  は周辺分布、  $p(\mathbf{y}_i|\mathbf{x}_j)$  は条件付

き確率、 $p(\mathbf{x}_j, \mathbf{y}_i)$  は統合確率 (同時確率) を表しています。最終的に得られた量は、分布  $p(\mathbf{x}, \mathbf{y})$  の下での  $-\log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  の期待値となることが分かります。これを次のように表現することにします。

$$\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})} \left[ -\log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) \right] = \begin{cases} \sum_i \sum_j \left[ -p(\mathbf{x}_j, \mathbf{y}_i) \log p_{\text{model}}(\mathbf{y}_i|\mathbf{x}_j; \boldsymbol{\theta}) \right] & \text{(離散確率分布)} \\ \int \int \left[ -p(\mathbf{x}, \mathbf{y}) \log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) \right] d\mathbf{x} d\mathbf{y} & \text{(連続確率分布)} \end{cases}$$

連続分布の場合は、 $\sum$  が  $\int$  になるだけです。この量は厳密に言えば「交差エントロピーの統合確率分布に関する期待値」なのですが、一般にはこの量のことを単に「交差エントロピー」と呼ぶことになっています。

※ ここでは少し教科書 “Deep Learning (2016)” とは異なる記法を用いているので、その理由についてコメントしておきます。教科書の式は、全体的に確率分布の引数が省略されてしまっていて、

$$\mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [\dots], \quad \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\dots], \quad \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})} [\dots]$$

は全て、

$$\mathbb{E}_{\mathbf{y} \sim p} [\dots], \quad \mathbb{E}_{\mathbf{x} \sim p} [\dots], \quad \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p} [\dots],$$

と表記されています。つまり、統合確率、周辺確率、条件付き確率をハッキリとは書き分けておらず、慣れている人からしたらスッキリして見やすいですし意味も分かるのですが、初めて勉強する人には分かりづらい書き方になっていると思います。この解説プリントでは、きちんとこれらを区別して書くことにしました。

## ♣ 損失関数としての交差エントロピー (2)

ここまでで、交差エントロピー

$$\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})} \left[ -\log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) \right]$$

が、真の確率分布  $p(\mathbf{y}|\mathbf{x})$  とモデル  $p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  の離れ具合を表す量となっていることを説明しました。ただし、この量をそのまま損失関数として用いることはできません。なぜなら、上の量を具体的に計算するためには、真の確率分布  $p(\mathbf{y}|\mathbf{x})$  やそれに対応する統合分布  $p(\mathbf{x}, \mathbf{y})$  を知っている必要があるからです。しかし、そもそもそれらを知らないからモデルを苦労して作るわけで、 $p(\mathbf{y}|\mathbf{x})$  や  $p(\mathbf{x}, \mathbf{y})$  は知っているはずがありません (知っているならモデルを作る必要がありません)。

では、どうしたらよいのでしょうか？それには、既に実際の測定などにより得られた手元にあるデータに基づいて  $p(\mathbf{x}, \mathbf{y})$  を近似することになります。この近似には様々な手法がありますが、1つの方法として、データから得られる「**経験分布 (empirical distribution)**」  $\hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})$  を使用する方法があります。つまり、 $\hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})$  が真の確率分布  $p(\mathbf{x}, \mathbf{y})$  に近いはずと考え、 $p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  を  $\hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})$  に近づけることで学習を行います。つまり、この2つの分布の交差エントロピー (の期待値)<sup>12</sup>

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} \left[ -\log p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) \right]$$

が、実際には損失関数として用いられます<sup>13</sup>。現在の多くのニューラルネットワークにおいて、この交差エントロピー誤差が損失関数として用いられています。

ここで出てきた経験関数は、実際に観測されたデータから得られる頻度割合の分布です。例えば、30人のクラスで1週間に読む本の数  $x$  を調べたところ、1冊が15人、2冊が10人、3冊が5人となったとしましょう。経験分布は、それぞれの値の全データ数に対する割合として与えられます。つまり、

$$\hat{p}_{\text{data}}(x=1) = \frac{15}{30} = \frac{1}{2}, \quad \hat{p}_{\text{data}}(x=2) = \frac{10}{30} = \frac{1}{3}, \quad \hat{p}_{\text{data}}(x=3) = \frac{5}{30} = \frac{1}{6}$$

となります。このようにすると、この経験分布  $\hat{p}_{\text{data}}(x)$  は確率の基本ルール「全て足したら1になる」

$$\sum_i \hat{p}_{\text{data}}(x_i) = 1$$

を満たすようになり、確率として扱うことができるようになります。教科書“Deep Learning (2016)”では、経験分布をもう少し数学的に次のように定義しています。

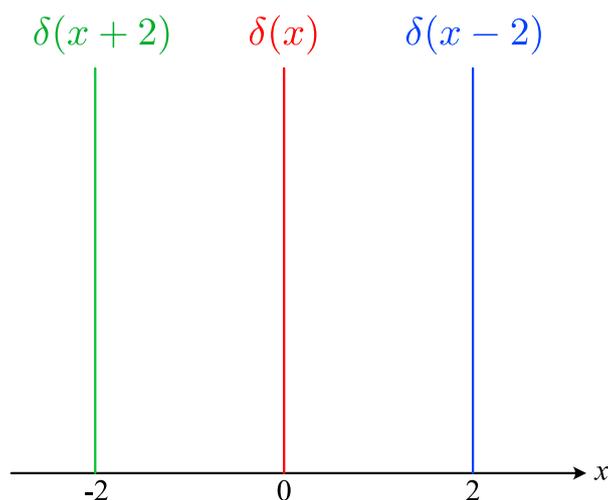
$$\hat{p}_{\text{data}}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} - \mathbf{x}^{(i)})$$

ここで出てきた  $\delta(\mathbf{x})$  は「デルタ関数」と呼ばれる特殊な関数です<sup>14</sup>。  $\delta(x - \mu)$  のグラフ (のイメージ) を描くと次のようになります。

<sup>12</sup>教科書“Deep Learning (2016)”の(6.12)式です。

<sup>13</sup>一般に経験分布  $\hat{p}_{\text{data}}$  を使って得られる期待値のことを「経験リスク (empirical risk)」といいます。

<sup>14</sup>「超関数」と呼ばれる関数の1つです



つまり、 $\delta(\dots)$  は、かっこの中が0となる時にピークを持ち、それ以外の値の時には0となるような関数です。ただしこれはあくまでイメージで、むしろ次の式をデルタ関数の定義として理解する方が正確です。

$$\sum_i f(\mathbf{x}_i) \delta(\mathbf{x}_i - \boldsymbol{\mu}) = f(\boldsymbol{\mu}) \quad (\mathbf{x} \text{ は離散的な変数})$$

$$\int f(\mathbf{x}) \delta(\mathbf{x} - \boldsymbol{\mu}) d\mathbf{x} = f(\boldsymbol{\mu}) \quad (\mathbf{x} \text{ は連続的な変数})$$

$\sum$  や  $\int$  によって、変数  $\mathbf{x}$  は様々な値を取ります。しかし、デルタ関数は基本的には0なので、ほとんどの  $\mathbf{x}$  の値に対して  $f(\mathbf{x})\delta(\mathbf{x} - \boldsymbol{\mu})$  は0となります。ただ唯一の例外として、デルタ関数の引数が0となる時 ( $\mathbf{x} = \boldsymbol{\mu}$  の時)、デルタ関数が値を持ち、その時の関数の値  $f(\mathbf{x} = \boldsymbol{\mu})$  が取り出されます。このように、総和 (積分) の中から特定の値を取り出すためにデルタ関数は用いられます。つまり、デルタ関数は単独でボンっと出てくることはなく、基本的には総和記号や積分の中でのみ登場する関数となっています。

2つの確率変数  $\mathbf{x}$  と  $\mathbf{y}$  の統合確率を表す経験分布も同じようにデルタ関数を使って表すことができます。

$$\hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} - \mathbf{x}^{(i)}) \delta(\mathbf{y} - \mathbf{y}^{(i)})$$

ここで、

$$(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \quad (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \quad \dots \quad (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$$

が  $m$  個の手持ちの訓練データを表します。この経験分布を先程の損失関数としての交差エントロピーの式に代入して、計算を進めてみましょう。

$$\begin{aligned}
J(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} \left[ -\log p_{\text{model}}(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}) \right] \\
&= \frac{1}{m} \sum_i \sum_j \left[ - \left( \sum_k \delta(\mathbf{x}_i - \mathbf{x}^{(k)}) \delta(\mathbf{y}_j - \mathbf{y}^{(k)}) \right) \log p_{\text{model}}(\mathbf{y}_j | \mathbf{x}_i; \boldsymbol{\theta}) \right] \\
&= \frac{1}{m} \sum_k \left[ -\log p_{\text{model}}(\mathbf{y}^{(k)} | \mathbf{x}^{(k)}; \boldsymbol{\theta}) \right]
\end{aligned}$$

このあたりはかなりややこしい計算になっているので、初めての方は何となく分かれば十分です。2行目から3行目に移る際には次のような操作を行なっています。期待値を取るためのループ  $\sum_i \sum_j$  によって、デルタ関数の引数に含まれる  $\mathbf{x}_i$  と  $\mathbf{y}_j$  が更新されてきます<sup>15</sup>。しかし、デルタ関数は基本的に0なので、ほとんどの  $\mathbf{x}_i$  と  $\mathbf{y}_j$  の値に対して、このループの中の関数の値は0となります。ただ唯一の例外として、 $(\mathbf{x}_i, \mathbf{y}_j) = (\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$  となる場合のみ、デルタ関数は値を持ち、先ほどと同様に  $\log p_{\text{model}}(\mathbf{y}^{(k)} | \mathbf{x}^{(k)}; \boldsymbol{\theta})$  という値が取り出されています。

色々ややこしい計算をしてきましたが、実用的には次の量 (交差エントロピー誤差) を損失関数として用いれば良いことになります<sup>16</sup>。

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_i \left[ -\log p_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \right]$$

$\frac{1}{m}$  は単なる定数なので、要するに

$$J(\boldsymbol{\theta}) = - \sum_i \log p_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$$

が最小となるように学習を行えば良いことになります。お気づきでしょうか、これは前回の「条件付き対数尤度と平均二乗誤差」でやった最尤法 (最尤推定) に他なりません。最尤法では、

$$\sum_i \log p_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$$

という**条件付き対数尤度**を最大化することで、最も適切なモデルのパラメータ  $\boldsymbol{\theta}_{\text{ML}}$  (最尤推定量) を決定しました。交差エントロピーは、この条件付き対数尤度にマイナスをつけた、いわば**「負の対数尤度」**になっているのです。交差エントロピーが現在多くのニューラルネットワークで損失関数として用いられている背景には、この量を最小化することが最尤推定に対応しており、統計学的にも最もらしい方法になっているという事情があるのです<sup>17</sup>。

<sup>15</sup>  $\sum_k$  は、経験分布を構成するためのループなので、 $i$  と  $j$  とは異なる添字 (カウンタ変数) を用いています。

<sup>16</sup> 添え字はなんでも良いので、ここでは  $k$  でなく  $i$  を使っています。

<sup>17</sup> 学習速度の面からも、交差エントロピーは非常に効率のよい損失関数となります。

## ♣ 損失関数としての平均二乗誤差と平均絶対誤差

ここまでで、損失関数としては交差エントロピー誤差

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} \left[ -\log p_{\text{model}}(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}) \right]$$

を用いれば良いことが分かりました。この交差エントロピー誤差では、特に  $p_{\text{model}}$  に対して具体的な分布を仮定しませんでした。しかし、ある特定の分布の形を仮定すると、この交差エントロピー誤差をさらに書き換えることができます。

1つ目として、モデルの確率分布が次のような Gauss 分布で与えられる場合を考えましょう<sup>18</sup>。

$$p_{\text{model}}(y | x; \boldsymbol{\theta}) = \mathcal{N}(y; f(x; \boldsymbol{\theta}), 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}[y - f(x; \boldsymbol{\theta})]^2\right)$$

分散は何でもよいのですが、ここでは簡単のために1としています。平均は入力値  $x$  とパラメータ  $\boldsymbol{\theta}$  で表現される関数としています。例えば、最初にあげた「走行速度」と「停車距離」のデータに対しては、おおまかに1次関数的な傾向が見て取れるので、

$$f(x; \boldsymbol{\theta}) = \theta_0 + \theta_1 x$$

とするのが良いでしょう。もちろん1次関数でなくて

$$f(x; \boldsymbol{\theta}) = \theta_0 + \theta_1 x + \theta_2 x^2$$

のような2次関数でも何でも大丈夫です。モデル分布がこのように与えられた時の交差エントロピー誤差を計算すると、

$$\begin{aligned} J(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} \left[ -\log \left( \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}[y - f(x; \boldsymbol{\theta})]^2\right) \right) \right] \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} \left[ -\log \frac{1}{\sqrt{2\pi}} + \frac{1}{2}[y - f(x; \boldsymbol{\theta})]^2 \right] = -\log \frac{1}{\sqrt{2\pi}} + \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} [y - f(x; \boldsymbol{\theta})]^2 \end{aligned}$$

となります。何回か出てきた対数  $\log$  の性質や、期待値の足し算の性質を使っています。また、定数の期待値に関しては、確率の合計が1となることに注意して、

$$\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} \left[ -\log \frac{1}{\sqrt{2\pi}} \right] = -\log \frac{1}{\sqrt{2\pi}} \sum_i \sum_j \hat{p}_{\text{data}}(\mathbf{x}_i, \mathbf{y}_j) = -\log \frac{1}{\sqrt{2\pi}}$$

となることを用いています。ただ、この部分は単なる定数なので、損失関数として使う際には

<sup>18</sup>簡単のためにここでは1変数の場合を考えます。

不要となります。つまり、モデルの分布が Gauss 分布の場合には、次の量を損失関数として使えば良いことが分かります<sup>19</sup>。

$$J(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E}_{x,y \sim \hat{p}_{\text{data}}(x,y)} [y - f(x; \boldsymbol{\theta})]^2 = \frac{1}{2m} \sum_{i=1}^m [y^{(i)} - f(x^{(i)}; \boldsymbol{\theta})]^2$$

この量は、前回の「条件付き対数尤度と平均二乗誤差」にて説明した、平均二乗誤差 (MSE: mean squared error)

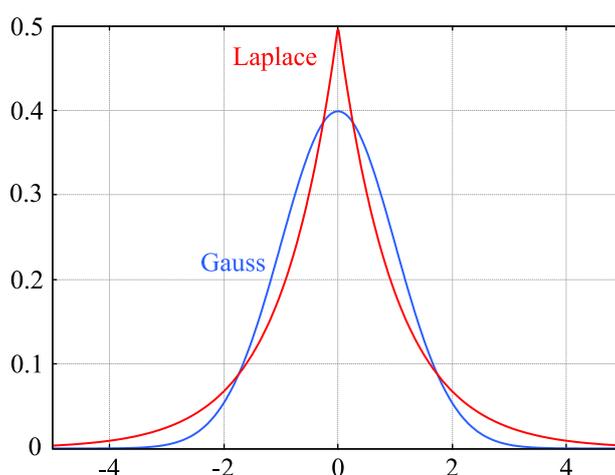
$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m [y^{(i)} - f(x^{(i)}; \boldsymbol{\theta})]^2$$

と定数倍  $\frac{1}{2}$  の違いを除いて全く同じであることが分かります。つまり、モデルの分布が Gauss 分布であると仮定した場合には、損失関数として交差エントロピーを使うことは平均二乗誤差を最小化することと一致するのです。インターネットで検索すると「損失関数としては平均二乗誤差と交差エントロピー誤差の2つが広く使われおり、問題によって使い分ける」といった記述をよく見かけるため、両者は全く別の量に勘違いしてしまいがちですが、今見たように本質的には同じものなのです。

先ほどは、分布が Gauss 分布の場合を考えましたが、今度は次のような Laplace 分布で与えられる場合を考えましょう。

$$p_{\text{model}}(y|x; \boldsymbol{\theta}) = \text{Laplace}(y; f(x; \boldsymbol{\theta}), 1) = \frac{1}{2} e^{-|y - f(x; \boldsymbol{\theta})|}$$

Laplace 分布のグラフは次のようになります。比較のために、分散と平均の値が等しい Gauss 分布も描いています。



(平均 0、分散 1 の Gauss 分布と Laplace 分布)

<sup>19</sup>教科書 “Deep Learning (2016)” の (6.13) 式に対応する式です。

グラフをみると、Laplace 分布は Gauss 分布と比べて裾野が広い分布になっていることが分かります。ですから、データの値が散らばっている場合には、Gauss 分布を用いるより Laplace 分布を用いた方が精度の良いモデルになる傾向にあります。この時の交差エントロピーを計算すると、

$$\begin{aligned} J(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} \left[ -\log \left( \frac{1}{2} e^{-|y - f(x; \boldsymbol{\theta})|} \right) \right] \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} \left[ -\log \frac{1}{2} + |y - f(x; \boldsymbol{\theta})| \right] = -\log \frac{1}{2} + \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} |y - f(x; \boldsymbol{\theta})| \end{aligned}$$

となります。先ほど同様、 $-\log \frac{1}{2} = \log 2$  の部分は単なる定数なので、この項を落として

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\mathbf{x}, \mathbf{y})} |y - f(x; \boldsymbol{\theta})| = \frac{1}{m} \sum_{i=1}^m |y^{(i)} - f(x^{(i)}; \boldsymbol{\theta})|$$

を損失関数として用いれば良いこととなります。この量は、「**平均絶対誤差**」(MAE: Mean Absolute Error)と呼ばれます。この量も、あくまで交差エントロピー誤差の具体例に過ぎません。

以上のように、モデルの確率分布を指定すると、交差エントロピーから平均二乗誤差や平均絶対誤差が導かれることが分かりました。ただ実際には、学習の効率の面から、交差エントロピーをそのまま損失関数として用いることの方が一般的となっています。