

楽しい微分可能プログラミング

JAX/NumPyro で遊ぼう v0.5.1

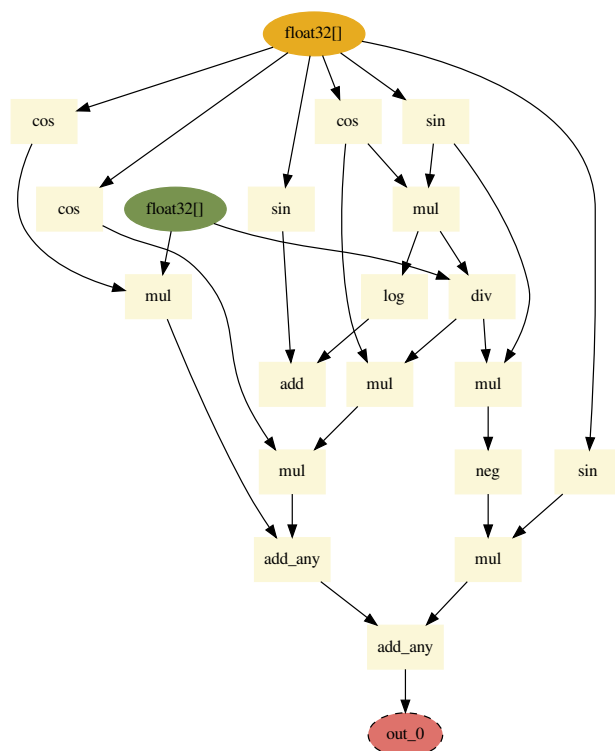


図1 $\partial_x(\log \cos x \sin x + \sin x)$ の計算グラフ

河原創

© 2020–2025 Hajime Kawahara

2025 年 6 月 24 日

目次

第 1 章	イントロダクション	3
1.1	自動微分の簡単な紹介	3
1.2	確率プログラミング言語	6
第 2 章	JAX の自動微分で遊ぶ	7
2.1	自動微分	7
2.2	for ループを自動微分する	14
2.3	JIT とクラス	19
第 3 章	微分可能プログラミングの基礎	21
3.1	方向微分と Jacobian Vector Product	21
3.2	Vector Jacobian Product	24
3.3	JAX の自動微分をカスタマイズする	26
3.4	Newton-Raphson 法で実装された陰関数の VJP を定義する	29
第 4 章	JAX プログラミングの様々なテクニック	33
4.1	XLA の再コンパイルを防ぐ	33
4.2	reverse-mode 微分と checkpoint	33
4.3	jax.numpy.array の範囲外アサイン	33
4.4	jax.numpy.array のインデックス操作	34
4.5	scan 中のインデクシング	35
第 5 章	JAX で最適化	37
5.1	jax.experimental	37
5.2	JAXopt	40
第 6 章	pytree/tree-math (TBD)	41
第 7 章	Numpyro でマルコフ鎖モンテカルロ・シミュレーションをする	42
7.1	曲線フィット	42
7.2	カスタマイズした自動微分で numpyro を動かす	52
第 8 章	NumPyro で楽しむガウス過程	55
8.1	ガウス過程	55

8.2	ガウス過程ノイズを含んだモデルフィット	64
第 9 章	NumPyro	66
9.1	scipy/JAX/NumPyro でサンプリング	66
第 10 章	orbax で pytree を保存する	68
10.1	save/restore	68
付録 A	Trouble shooting	71
A.1	JAX 編	71
付録 B	TIPS	72
B.1	JAX	72
B.2	Arviz	73
付録 C	Appendix	74
C.1	A. ガウス分布	74
参考文献		75

第 1 章

イントロダクション

JAX は google が開発中の自動微分と XLA(accelerated Linear Algebra) パッケージです [2]。つまり自動微分+線形代数パッケージで、なおかつ `numpy` と同じ形式で書けるようになっています。本稿はこの JAX と PPL として特に `numpyro`、また最適化ツール `JAXopt` を使って遊びたいというための文書です。普通であれば「JAX/numpyro で学ぶ機械学習」のようなタイトルになると思うのですが、むしろ私は JAX/numpyro で遊びたいのです。というのも、毎年、年末年始になると時間ができるので新しい言語やパッケージで遊びたいのです。2021 年はコロナ下ということもあり JAX で[太陽系外惑星のスペクトル推定コード](#)を書いてみようというのをやってみて[査読論文](#)にしました。最近[続編](#)もだしました。さて、最近 JAX をはじめとする自動微分を実装したプログラミングを微分可能プログラミング (Differentiable Programming) と呼び、新たなパラダイムとみなすことが多くなってきました。そこでえ 2024 年末にタイトルを変更し「微分可能プログラミング」が前面に出るようにしました。これには”The Elements of Differentiable Programming” の登場 [1] が大きかったように思います。随時更新する予定です*¹。

本稿で用いられているコードは[github](#)にいたのでご利用ください。

1.1 自動微分の簡単な紹介

勾配ベースの最適化や HMC-NUTS、変分推定にはモデルを最適化あるいは推定したいパラメタで微分できる必要がある。数値的に微分値を求めるにはおおまかに 4 つの方法がある。一つ目は手で微分 (manual differentiation) をし、結果をコーディングすることである。これはモデルが複雑になってくると破綻しがちであり、またフレキシビリティの観点からもモデルの継続的な改良を妨げる。次に `mathematica` 等による symbolic differentiation を用いて手で微分する代わりに微分結果を得てコーディングすることが考えられる。これは `mathematica` 等を用いたことのある方ならばわかると思うが、モデルが複雑になってくると膨大な項数の結果が排出されるので、やはりフレキシビリティの観点からは難点がある。次に、数値微分を行うことが考えられる。数値微分はモデルが複雑になってくるとエラーがたまりやすい。そこで機械学習分野などで用いられているのが JAX でも用いられている自動微分である。コーディングできるモデルは通常、様々な微分の既知である関数（ここでは要素関数と呼ぼう）の加減乗除の組み合わせでできている。そこで、自動微分では、

*¹ 質問・コメント・お問い合わせはdivrot@gmail.comまでお寄せください。なお筆者はソプラノサクソプレイヤー兼野生の[天文学の研究者](#)です。

連鎖則が成り立つように各要素関数

$$x \mapsto f(x) \quad (1.1)$$

を関数のアウトプットとその微分値の情報を出力するように拡張する。また微分の加乗除の演算規則がなりたつように演算を定義する。自動微分の実装法としては JAX では Jacobian Vector Product (JVP) を用いたものが使用されている。しかしここではイントロダクションとしてより説明が容易な、双対数を用いた方式を使用して自動微分を解説しよう。

双対数 (dual number)、 $z \in k[\epsilon]/\langle \epsilon^2 \rangle$ ^{*2} は $a, b \in \mathbb{R}$ にたいし、

$$z = a + b\epsilon \quad (1.2)$$

$$\epsilon^2 = 0 \quad (1.3)$$

となる数である。複素数は $i^2 = -1$ であったのが $\epsilon^2 = 0$ となったと考えればよい。変数 x の拡張として実部 a に x を、非実部 b に x' を割り当てると、 $z = f + f'\epsilon, w = g + g'\epsilon$ に対し、可算・乗数・除算はそれぞれ

$$z + w = (f + g) + (f' + g')\epsilon \quad (1.4)$$

$$zw = fg + (f'g + fg')\epsilon \quad (1.5)$$

$$z/w = \frac{f}{g} + \frac{f'g - fg'}{g^2}\epsilon \quad (1.6)$$

となる。これは実部に関して通常の、また非実部に関して微分の加乗除則と同じであることを示している。次に連鎖則を実現するには関数 $x \mapsto F(x)$ を

$$x + x'\epsilon \mapsto F(x) + F'(x)x'\epsilon \quad (1.7)$$

と拡張すればよい。ここで $F(x)$ を元の関数、双対数をもつ拡張された関数を $\hat{F}(x + x'\epsilon)$ と分けて表記する。すると式 (1.7) は

$$\hat{F}(x + x'\epsilon) = F(x) + F'(x)x'\epsilon \quad (1.8)$$

と書ける。さて連鎖則 $G(F(x))$ の拡張は

$$\hat{G}(\hat{F}(x + x'\epsilon)) = \hat{G}(F(x) + F'(x)x'\epsilon) \quad (1.9)$$

$$= G(F(x)) + G'(F(x))F'(x)x'\epsilon \quad (1.10)$$

となり、実部は通常の合成 $(G(F(x)))$ が、非実部は連鎖則 $G'(F(x))F'(x)x' = \frac{dG}{dF} \frac{dF}{dx} x'$ が実現されているのが確認できる。

1.1.1 ミニマム自動微分

というわけで小さな自動微分を Python で実装してみよう。ここで解きたい問題を

$$F(x) = \log(\cos x \sin x) + \sin x \quad (1.11)$$

^{*2} 多項式環の商環のことをしめしている。詳しくはたとえば雪江明彦「環と体とガロア理論」などを参照。

の $x = 1$ で微分 $F'(1)$ とする。さらにその合成関数

$$\begin{aligned} G(x) &= F(F(x)) \\ &= \log(\cos(\log(\cos x \sin x) + \sin x) \sin(\log(\cos x \sin x) + \sin x)) + \sin(\log(\cos x \sin x) + \sin x) \end{aligned} \quad (1.12)$$

の $x = 1$ での微分 $G'(1)$ も求めたい。式 (1.12) を手で微分するのは大変であるし、symbolic differentiation の結果も長大である。しかし自動微分ならばかなりすっきりしたコードになる。

まず関数 (1.11) に存在する演算は可算と乗算であるので、式 (1.4)、(1.5) に対応した双対数の可乗算を定義する。

```
1 def mul(x, y):
2     a, b = x
3     c, d = y
4     return a*c, a*d + b*c
5
6 def add(x, y):
7     a, b = x
8     c, d = y
9     return a + c, b + d
```

ここに x, y はそれぞれ双対数である。次に関数 (1.11) に存在する要素関数は $\sin x$ 、 $\cos x$ 、 $\log x$ である。式 (1.8) を実装するには、双対数の実部と非実部をペア (a, b) の形式であらわすと、関数 $F(x)$ に対し

$$(x, dx) \rightarrow (F(x), F'(x)dx) \quad (1.13)$$

となるように入出力を与えればよいことがわかる。そこで

```
1 import numpy as np
2 def cos(x):
3     a, b = x
4     return np.cos(a), - np.sin(a)*b
5
6 def sin(x):
7     a, b = x
8     return np.sin(a), np.cos(a)*b
9
10 def log(x):
11     a, b = x
12     return np.log(a), b/a
```

これで終わりである。あとは

```
1 f = lambda x: add(log(mul(cos(x), sin(x))), sin(x))
```

PPL	backend
NumPyro	JAX
Tensorflow Probability	Tensorflow
Pyro	PyTorch
PyMC	Aesara/JAX
BlackJAX*	JAX

表1.1 代表的な PPL とそのバックエンド。*: BlackJAX は厳密には PPL ではなくサンプラーである。さまざまな PPL 上で利用可能。

```

2 df = lambda x: f([x,1.0])
3 df(1.0) #->(0.05324076815279066, -0.37501280285243144)

```

とすれば $(F(1), F'(1))$ が求まる。 $G'(x)$ も簡単に、

```

1 g = lambda x: f(f(x))
2 dg = lambda x: g([x,1.0])
3 dg(1.0) #-> (-2.8816056725768977, -7.391555094461485)

```

とするだけである。以上の例からわかるように自動微分の実際の計算には代数計算を用いているため数値微分に比べて誤差の蓄積は少ない。またフレキシブルにコーディングが可能である。

ここではミニマムな自動微分の説明を行った。第3章ではもう少し踏み込んだ説明を行うので興味のある方は参照されたい。

1.2 確率プログラミング言語

確率プログラミング言語 (Probabilistic Programming Language; PPL) は確率モデルを扱うためのプログラミング言語のことである。様々なものが存在するが、ほとんどのものはバックエンドとして自動微分パッケージを用いている。表1.1に Python で動く代表的な PPL とそのバックエンドを示す。

本文書では JAX をバックエンドとする NumPyro を用いる。

第 2 章

JAX の自動微分で遊ぶ

JAX [2] は tensorflow や PyTorch 等と同様に自動微分 [7, 3] を利用することができる。またその構文は numpy と形式的にほぼ対応している。また、後の章では JAX を用いた最適化や MCMC などを行う。そのためまず JAX で様々な関数を作成できるようにしよう。本章ではまず JAX で一通り自動微分を使用してみたから、第3.3章の自動微分のカスタマイズの説明の際に、自動微分の原理の説明を行う構成である。逆に言うと自分で自動微分を定義する必要に迫られない限り原理の理解は必要ないところが JAX の素晴らしいところでもあるといえるかもしれない。

2.1 自動微分

JAX の自動微分で遊んでみよう。

2.1.1 一変数の自動微分

$$f(x) = |x| + \sin x \quad (2.1)$$

を自動微分してみよう。まず `jax.numpy` をインポートする^{*1}。

```
1 import jax.numpy as jnp
```

`jax.numpy` の関数を利用して

```
1 def f(x):  
2     return jnp.abs(x) + jnp.sin(x)
```

もしくはラムダ関数を利用して

```
1 f = lambda x: jnp.abs(x) + jnp.sin(x)
```

^{*1} `jax.numpy` は `numpy` と併用することが多く、`jax.numpy` を `np` で `numpy` を `onp` とする流派、`jax.numpy` を `jnp` で `numpy` を `np` とする流派がいるようである。本稿では後者を採用する。

のように関数 $f(x)$ を定義する。

自動微分は `grad` を用いる。

```
1 from jax import grad
```

微分を

```
1 df=grad(f)
```

とすると、例えば、 $f'(1.0)$ は

```
1 df(1.) # → DeviceArray(1.5403023, dtype=float32)
```

のようにして求まる^{*2}。二階微分も同様に

```
1 df=grad(df)
2 ddf(1.0) # → DeviceArray(-0.841471, dtype=float32)
```

のように求まる。

ベクトルマッピング

上では微分の引数として、スカラーを代入していた。以下のようなベクトルを代入するとエラーとなる。

```
1 xv=jnp.linspace(-10,10,100)
2 df(xv) # → エラー！
```

これには `jax.vmap` を用いてマッピングをすることで可能となる。

```
1 from jax import vmap
2 dfv=vmap(grad(f))
3 dfv(xarr) # → jnp.array で返ってくる。
```

これで図が描けるようになった。

```
1 import matplotlib.pyplot as plt
2 plt.plot(xarr,f(xarr),".",label="$f(x)=|x|+\sin\{x\}$")
3 plt.plot(xarr,dfv(xarr),".",label="f'(x)")
4 plt.xlabel("x")
5 plt.legend()
```

^{*2} `df(1)` だとエラーがでる。これは引数に `int` を取れないためである。

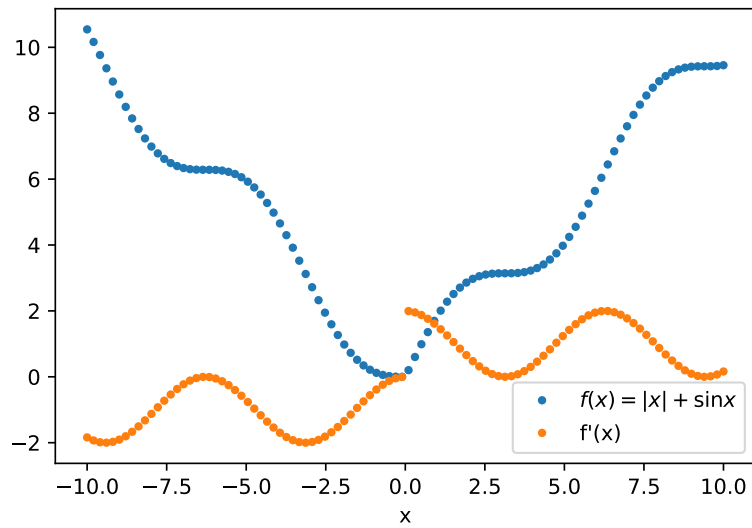


図2.1

2.1.2 多変数の自動微分

次に $\mathbf{r} = (x, y, z)^\top$ を変数とした多変数関数

$$g(\mathbf{r}) = \sqrt{x^2 + 2y^2 + 3z^2} \quad (2.2)$$

の自動微分を行おう。関数定義は

```
1 g = lambda x,y,z: jnp.sqrt(x**2+2*y**2+3*z**2)
```

である。偏微分

$$\partial_x g(\mathbf{r}) = \frac{\partial g(\mathbf{r})}{\partial x} \quad (2.3)$$

として $\mathbf{r}_1 = (1, 1, 1)^\top$ での $\partial_x g(\mathbf{r}_1)$ は

```
1 dgdx=grad(g, argnums=0)
2 dgdx(1.,1.,1.) # -> DeviceArray(0.40824828, dtype=float32)
```

ここで `argnums` が微分する変数を指定している。したがって

$$\partial_y g(\mathbf{r}) = \frac{\partial g(\mathbf{r})}{\partial y} \quad (2.4)$$

として $\mathbf{r}_1 = (1, 1, 1)^\top$ での $\partial_y g(\mathbf{r}_1)$ は

```

1 dgdy=grad(g, argnums=1)
2 dgdy(1.,1.,1.) # → DeviceArray(0.81649655, dtype=float32)

```

となる。

もしくは

$$\nabla g(\mathbf{r}) = \frac{\partial g(\mathbf{r})}{\partial \mathbf{r}} \quad (2.5)$$

は、

```

1 dgdr=grad(g, argnums=(0,1,2))
2 dgdr(1.,1.,1.)
3 # → (DeviceArray(0.40824828, dtype=float32),
4 # DeviceArray(0.81649655, dtype=float32),
5 # DeviceArray(1.2247448, dtype=float32))

```

となる。

マッピング

さて、多変数の場合も、バッチ入力する際、そのまま入力するとエラーとなってしまう。

```

1 xv=jnp.linspace(-10,10,100)
2 dgdr(xarr,1.,1.) # → エラー！

```

そこでやはり `vmap` を利用する。この場合、どの変数をマッピングするか指定する。

```

1 dgdrv=vmap(grad(g, argnums=(0,1,2)), (0, None, None), 0)

```

(0, None, None) の部分で指定している。

```

1 dd=dgdrv(xv,1.,1.)
2 plt.plot(xv,dd[0],label="$\partial_x g(x,1,1)$")
3 plt.plot(xv,dd[1],label="$\partial_y g(x,1,1)$")
4 plt.plot(xv,dd[2],label="$\partial_z g(x,1,1)$")
5 plt.xlabel("x")
6 plt.legend()

```

で、[図2.2](#)を描くことができる。

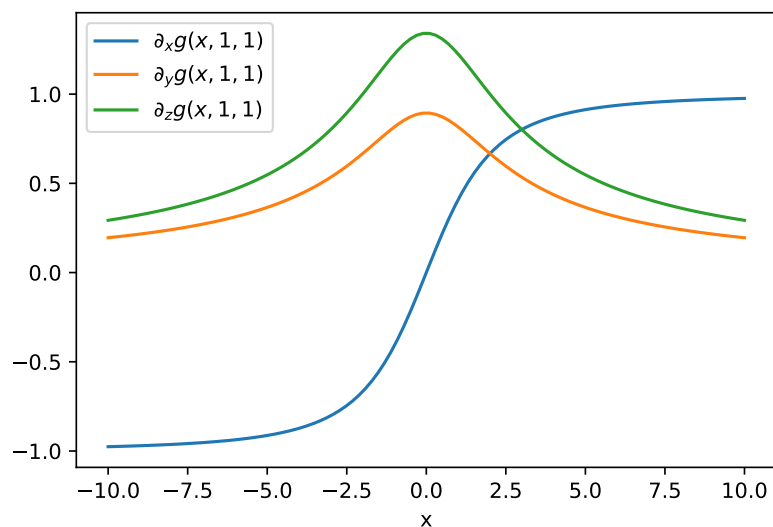


図2.2

2.1.3 入力がベクトルの場合

入力がベクトルの場合も考えよう。例えばとても簡単な

$$f(\boldsymbol{x}, \boldsymbol{y}) = \boldsymbol{x} \cdot \boldsymbol{y} \quad (2.6)$$

の場合を考えよう。

```
1 def f(x,y):
2     return jnp.dot(x,y)
```

この場合でも入力がスカラーだったときと同様に自動微分を計算できる。例えば $\partial_{\boldsymbol{x}} f(\boldsymbol{x}, \boldsymbol{y})$ は

```
1 x=jnp.linspace(0,1,11)
2 y=jnp.linspace(1,2,11)
3 f(x,y) # → DeviceArray(9.35, dtype=float32)
4 grad(f,argnums=0)(x,y) # → DeviceArray([1., 1.1,...,2.], dtype=float32)
```

のように計算できる。

`vmap` を用いれば複数の入力ベクトルに対し、マッピングできる。例えば、5つの \boldsymbol{x} を `vmap` で同時に処理してみよう。

```
1 X=jnp.linspace(0,1,55).reshape(5,11)
2 vmap(grad(f,argnums=0),(0,None),0)(X,y) # → DeviceArray([[1. ...])
```

ここでも `vmap` の `(0, None)` に注意しよう。5 つの (x, y) に対して同時に処理するには

```
1 X=jnp.linspace(0,1,55).reshape(5,11)
2 Y=jnp.linspace(1,2,55).reshape(5,11)
3 vmap(grad(f,argnums=0),(0,0),0)(X,Y) # → DeviceArray([[1. ...])
```

のように、`vmap` の引数を `(0,0)` にすれば良いこともわかる。

2.1.4 練習

微分が含まれた関数の表示

経済学で用いられる Utility function $u(x)$ を考えよう。Utility function は上に凸 (concave function) であるとして Absolute Risk Aversion (ARA) と Relative Risk Aversion (RRA)

$$\text{ARA}(x) \equiv -\frac{u''(x)}{u'(x)} \quad (2.7)$$

$$\text{RRA}(x) \equiv -\frac{u''(x)}{u'(x)}x \quad (2.8)$$

をプロットしてみよう。以下では、例として $u(x) = \log(1+x)$ としている。

```
1 #sample utility function
2 def u(x):
3     return jnp.log(1.0+x)
```

`vmap` を用いてマッピングした後に演算することで `jax numpy array` 形式を引数として受け付けるようになる。

```
1 def ARA(xarr):
2     den=vmap(grad(grad(u)))
3     bun=vmap(grad(u))
4     return -den(xarr)/bun(xarr)
```

```
1 def RRA(xarr):
2     den=vmap(grad(grad(u)))
3     bun=vmap(grad(u))
4     return -den(xarr)/bun(xarr)*xarr
```

プロットは

```
1 x_=jnp.linspace(0,3,100)
2 plt.plot(x_,u(x_),label="$u(x)$")
3 plt.plot(x_,ARA(x_)
4 ,label="ARA$(x) \equiv -u^{\prime\prime}(x)/u'(x)$",ls="dashed")
```

```

5 plt.plot(x_, RRA(x_))
6 ,label="RRA$(x) \equiv -u^{\prime\prime}(x)/u^{\prime}(x)x$",ls="dotted")
7 plt.legend()
8 plt.xlabel("$x$")
9 plt.xlim(0.,3.0)
10 plt.ylim(0,3)

```

のようにすればよい。図2.3に結果を示す。

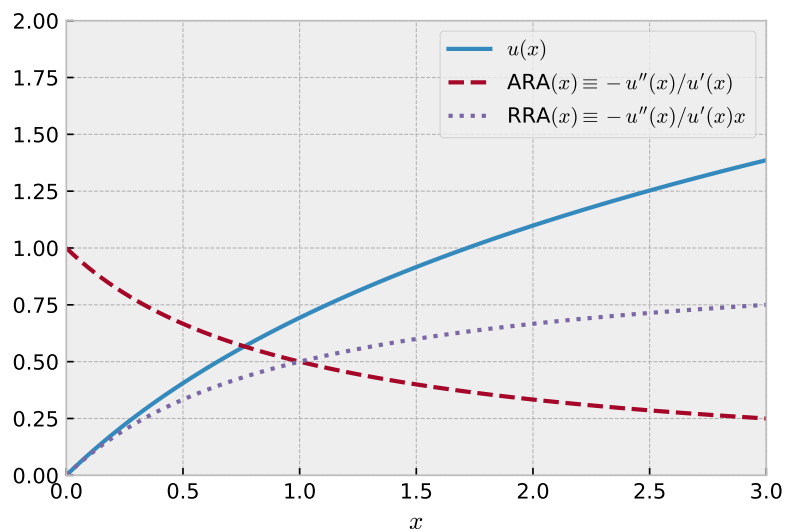


図2.3

コントロールパラメタを入れる

$$u(x, A) = \log(A + x) \quad (2.9)$$

のようにコントロールパラメタ A を入れた場合についても例示しよう。この場合、

```

1 def u(x,A):
2     return jnp.log(A+x)

```

を用いることになる。さてここでは更に ARA や RRA の微分を図示しよう。すなわち

$$\partial_x \text{ARA}(x) = \frac{\partial}{\partial x} \left(-\frac{u''(x)}{u'(x)} \right) \quad (2.10)$$

$$\partial_x \text{RRA}(x) \equiv \frac{\partial}{\partial x} \left(-\frac{u''(x)}{u'(x)} x \right) \quad (2.11)$$

をさまざまな A で図示したい。これには微分があくまで第 0 argument について行われていることに注意し、vmap は最後にこれまた第 0 argument についてマッピングすることに注意すると

```
1 def dARA(xarr,A):
2     den=grad(grad(u,argnums=(0)),argnums=(0))
3     bun=grad(u,argnums=(0))
4     g=lambda x,A:-den(x,A)/bun(x,A)
5     h=vmap(grad(g,argnums=(0)), (0, None), 0)
6     return h(xarr,A)
```

および

```
1 def dRRA(xarr,A):
2     den=grad(grad(u,argnums=(0)),argnums=(0))
3     bun=grad(u,argnums=(0))
4     g=lambda x,A:-den(x,A)/bun(x,A)*x
5     h=vmap(grad(g,argnums=(0)), (0, None), 0)
6     return h(xarr,A)
```

のように定義すれば良いことがわかる。図2.4は

```
1 Aarr=[-1.0,0.0,1.0]
2 fig=plt.figure(figsize=(15,3))
3 for i,Ain enumerate(Aarr):
4     ax=fig.add_subplot(1,3,i+1)
5     x_=jnp.linspace(-Ain+0.5,-Ain+3,100)
6     plt.title("A="+str(Ain))
7     ax.plot(x_,u(x_,Ain),label="$u(x)$")
8     ax.plot(x_,dARA(x_,Ain),label="$\partial_x ARA(x)$",ls="dashed")
9     ax.plot(x_,dRRA(x_,Ain),label="$\partial_x RRA(x)$",ls="dotted")
10    plt.legend()
11    plt.xlabel("$x$")
```

のように書ける。

2.2 for ループを自動微分する

さてこれまで比較的簡単な関数を自動微分してきた。この程度なら手で（もしくは sympy とか mathematica で）微分できるじゃないかといわれそうである。まあ別に究極的には手でできるのだが、いかにも面倒くさそうな場合として関数が for loop を含む場合^{*3}、自動微分を使えると便利そうである。

^{*3} 古代の言語 fortran ならば do ループのことである。

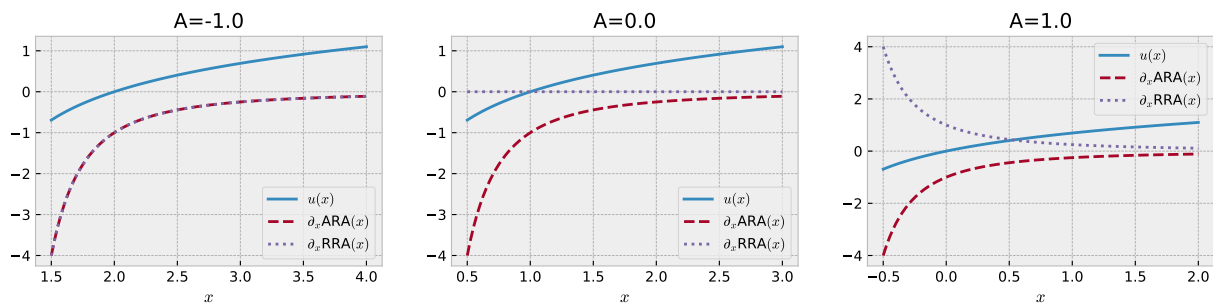


図2.4

jax では `fori_loop` というのがあるが、version 0.2.6 の時点ではこれを用いると自動微分が通らない。そこで用いるのが `jax.lax.scan` である。scan と等しい python コードは jax 公式チュートリアルによると

```

1 def scan(f, init, xs, length=None):
2     if xs is None:
3         xs = [None] * length
4     carry = init
5     ys = []
6     for x in xs:
7         carry, y = f(carry, x)
8         ys.append(y)
9     return carry, np.stack(ys)

```

である。これは for loop より一般的だが、for loop に使おうと思うとなんかちよっと頭がくらくらする。xs を None として、length を for loop の回数として、ほぼ同じコードなのだが以下のように書き換える。

```

1 def scan(f, x0, length=n):
2     step = [None] * n
3     x = x0
4     narr = []
5     for i in step:
6         x, null = f(x, i)
7         narr.append(null)
8     return x, np.stack(narr)

```

これは、特に使っていない null とか narr を無視すれば、ほとんど以下の for loop と同じである

```

1 def scan_eq(f, x0, n)
2     x=x0
3     for i in range(0,n):
4         x = f(x)

```



```
5     return x
```

というわけで、scan で for loop が実装できることが分かった。

さて、例として

$$f(x) = 1/(1+x) \quad (2.12)$$

を n 回適用させた関数 $g(x;n)$ を考えよう。例えば $n = 3$ の時、

$$g(x;3) = f(f(f(x))) = \frac{1}{1 + \frac{1}{1 + \frac{1}{1+x}}} \quad (2.13)$$

である。これを微分した関数

$$h(x;n) = \frac{\partial}{\partial x} g(x;n) \quad (2.14)$$

を計算する。

まず scan, grad を import する。

```
1 from jax.lax import scan
2 from jax import grad
```

関数 f を定義する。

```
1 def f(x,null):
2     x=1.0/(1.0+x)
3     return x,null
```

null は別に意味のない変数だが scan の形式に合わせるために無駄についている。次に scan で for loop 部分を実装する。今は $n = 3$ で試そう。

```
1 def g(x0):
2     x,null=scan(f,x0,None,3)
3     return x
```

そして自動微分を定義する。ここでベクトル化した xarr を代入するために vmap しておこう。jit というのは Just-In-Time コンパイラ（実行時コンパイラ）のことで、コンパイルするので高速になる。今回の例では jit をするのとしないので雲泥の差が出るので試されたし。

```
1 from jax import vmap
2 from jax import jit
3 h=jit(vmap(grad(g)))
4 xarr=jnp.linspace(0,1,100)
5 harr=h(xarr)
```

さて、答え合わせを行うために sympy でも計算しておこう。

```
1 import sympy as sp
2 x = sp.Symbol('x')
3 def f_sp(x):
4     return 1/(1+x)
5 sp.diff(f_sp(f_sp(f_sp(x))),x)
```

答えが、

$$h(x,3) = -\frac{1}{\left(1 + \frac{1}{1+\frac{1}{x+1}}\right)^2 \left(1 + \frac{1}{x+1}\right)^2 (x+1)^2} \quad (2.15)$$

と求まる。つまり比較すべき関数は

```
1 h_sp=lambda x:-(f_sp(f_sp(f_sp(x))))**2*(f_sp(f_sp(x)))**2*(f_sp(x))**2
```

となる。比較の図は、

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 plt.style.use('bmh')
4
5 fig=plt.figure()
6 ax=fig.add_subplot(211)
7 plt.plot(xarr,h_sp(xarr)-harr)
8 plt.ylabel("sympy-jax")
9 ax=fig.add_subplot(212)
10 plt.plot(xarr,harr)
11 plt.xlabel("$x$")
12 plt.ylabel("$\partial_x f(f(f(x)))$")
13 plt.show()
```

でかけ、結果は、図2.5である。このように直接関数定義した場合と 6,7 桁くらいの精度で一致しているのがわかる。

2.2.1 練習

scan のインプット配列で微分する

jax.lax.scan の他の使い方を練習してみよう。例えば $\mathbf{x} = (x_0, x_1, x_2, x_3)^\top$ に対して関数

$$f(\mathbf{x}) = \left(\left(\left(\left(x_0 + \frac{1}{x_0} \right) x_1 + \frac{1}{x_1} \right) x_2 + \frac{1}{x_2} \right) x_3 + \frac{1}{x_3} \right) \quad (2.16)$$

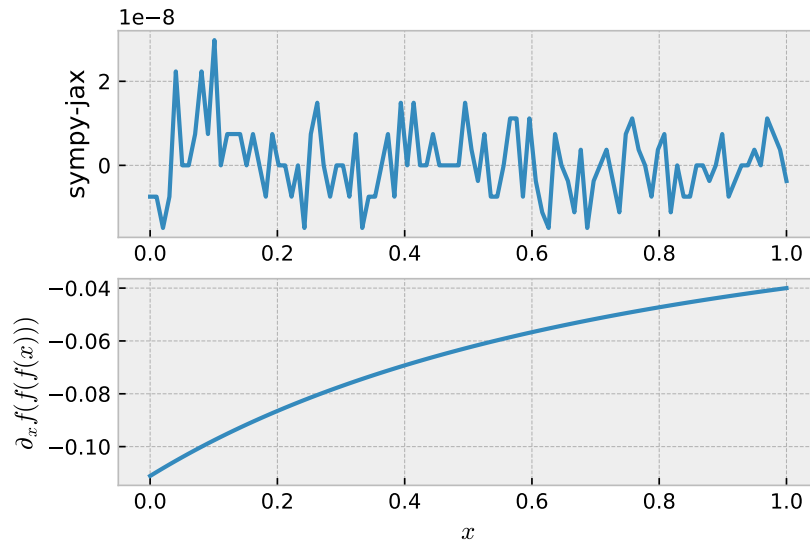


図2.5

を定義する。この関数の微分

$$\mathbf{h}(\mathbf{x}) \equiv \nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x_0}, \frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \frac{\partial f(\mathbf{x})}{\partial x_3} \right)^\top \quad (2.17)$$

の自動微分を実装してみよう。`scan(f, init, xs)` の `xs` を \mathbf{x} に用いればよいことが予想される。式 (2.16) は内側のかっこから再帰的に $(\cdot)x_i + 1/x_i$ を計算していけばよいことに気づく。そこで、まず

```
1 def g(y,x):
2     y=y*x+1.0/x
3     return y,None
```

とする。そして `scan` を用いて

```
1 @jit
2 def f(xs):
3     y0=1.0
4     y,null=scan(g,y0,xs)
5     return y
```

とすれば再帰的に $y_i = y_{i-1}x_i + 1/x_i$ を計算できることがわかる。この自動微分は単純に

```
1 h=grad(f)
```

とすればよい。これで、例えば、

```

1 xs=jnp.array([1.,2.,4.,5.])
2 h(xs) # → DeviceArray([ 0.      , 35.      , 22.1875, 18.21  ], dtype=float32)

```

と求まる。この形式はいくらでも要素数を増やせる利点がある。すなわち例えば、

```

1 h(jnp.array([1.,2.,3.,4.,5.,6.,7.,8.]))
2 # → DeviceArray([0., 35280., 29493.332, 23135.
3 # , 18662.559, 15572.979 , 13350.7705, 11682.194 ], dtype=float32)

```

といった感じである。

2.3 JIT とクラス

クラスを定義してクラス内の中においた関数に JIT をおこなえると便利である。この場合は、以下のように `functools.partial` を用いてデコレータをつけるとよい。

```

1 from functools import partial
2 class fcjax(object):
3     def __init__(self,c,n):
4         self.c=c
5         self.n=n
6
7     def f(self,x,null):
8         x=self.c/(self.c+x)
9         return x,null
10
11     @partial(jit, static_argnums=(0,))
12     def g(self,x0):
13         x,null=scan(f,x0,None,self.n)
14         return x
15
16     @partial(jit, static_argnums=(0,))
17     def h(self,x0):
18         hh=vmap(grad(self.g))
19         return hh(x0)

```

このようにクラスで書くことにより、例えば、関数の中を引数を変えて

```

1 for i in range(1,4):
2     fcc=fcjax(1.0,i)

```

```
3     plt.plot(fcc.h(xarr),label="n="+str(i))
4 plt.legend()
5 plt.show()
```

のように使うことができる。また、JIT で実行時にコンパイルされてしまうので、一度、`g` や `h` を呼び出してしまうと、その後に `fcc.n=4` のように更新しても、関数の方の `fcc.n` は変更されないことに注意が必要である

```
1 fcc=fcjax(1.0,1)
2 fcc.n=2
3 fcc.h(xarr) # この時点でコンパイルされ、n=2で計算される。
4 fcc.n=3
5 fcc.h(xarr) # すでにコンパイルされてしまっているのでやはり n=2のままで計算される
```

第 3 章

微分可能プログラミングの基礎

3.1 方向微分と Jacobian Vector Product

方向微分とは、多変数関数において、特定の方向に沿った変化率を評価するものである。与えられた点 ω において関数の変化率を、その点を基準とし任意のベクトル方向 \mathbf{v} に沿って計算する。関数 $f(\omega)$ の $\mathbf{v} \in \mathbb{R}^D$ 方向の方向微分は以下で定義される。

$$\partial f(\omega)[\mathbf{v}] \equiv \lim_{\delta \rightarrow 0} \frac{f(\omega + \delta \mathbf{v}) - f(\omega)}{\delta} \quad (3.1)$$

ω がベクトルの場合

ω がベクトルの時、つまり $\omega \in \mathbb{R}^D$ のとき、 $\delta \mathbf{v}$ が小さいとすると ($\delta|\mathbf{v}| \sim 0$)、Jacobian $\partial f(\omega) \in \mathbb{R}^{M \times D}$ をもちいて、

$$f(\omega + \delta \mathbf{v}) \sim f(\omega) + \partial f(\omega) \delta \mathbf{v} \quad (3.2)$$

となることから、

$$\partial f(\omega)[\mathbf{v}] = (\partial f) \mathbf{v} = \text{JVP}(\mathbf{f}, \mathbf{v}) \quad (3.3)$$

となることより、方向微分は Jacobian Vector Product (JVP) と呼ばれる。

ω が行列の場合

ω は行列でも良い。そこで ω を行列 W であるとする。

1. $f(W) = W\mathbf{x}$

行列による線形変換

$$f(W) = W\mathbf{x} \quad (3.4)$$

を考える。このとき、JVP は

$$\partial f(W)[V] \equiv \lim_{\delta \rightarrow 0} \frac{(W + \delta V)\mathbf{x} - W\mathbf{x}}{\delta} = V\mathbf{x} \quad (3.5)$$

となる。これは行列 W をベクトル化して式 (3.3) で JVP を計算したものと一致する。これを二次元の場合で例示する。

$$W = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \quad (3.6)$$

のとき、 $\omega = \text{Vec}(W) = (a, b, c, d)^T$ と同型変換すると、 $f(W) = Wx$ は

$$f(\omega) = X\omega, \quad X \equiv \begin{pmatrix} x & y & 0 & 0 \\ 0 & 0 & x & y \end{pmatrix} \quad (3.7)$$

となる。Jacobian は

$$\partial f(\omega) = \begin{pmatrix} \partial_a(ax+by) & \partial_b(ax+by) & \partial_c(ax+by) & \partial_d(ax+by) \\ \partial_a(cx+dy) & \partial_b(cx+dy) & \partial_c(cx+dy) & \partial_d(cx+dy) \end{pmatrix} = X \quad (3.8)$$

であるので、 V をベクトル化したものを $v = \text{vec}(V)$ とおくと、

$$\partial f(\omega)v = Xv \quad (3.9)$$

となっている。これは $Xv = Vx$ であることから

$$\partial f(W)[V] = Vx \quad (3.10)$$

であることが確認できる。

2. $f(W) = WX$

X も行列の場合を考える。つまり

$$f(W) = WX \quad (3.11)$$

の場合、JVP は同様に

$$\partial f(W)[V] \equiv \lim_{\delta \rightarrow 0} \frac{(W - \delta V)X - WX}{\delta} = VX \quad (3.12)$$

である。

3.1.1 JVP の連鎖則

Jacobian の連鎖則

$$\partial(g \circ f)(\omega) = \partial g(f(\omega))\partial f(\omega) \quad (3.13)$$

であることから

$$\partial(g \circ f)(\omega)[v] = \partial(g \circ f)(\omega)v = \partial g(f(\omega))\partial f(\omega)v = \partial g(f(\omega))[\partial f(\omega)[v]] \quad (3.14)$$

となる。つまり $g \circ f$ の方向微分は、図3.1に示すように、 g の評価点として $f(\omega)$ 、方向として f の v 方向の方向微分を用いた方向微分を計算すればよい。このような計算を繰り返すことにより JVP の合成関数の計算を順次的に行うことができる (図3.2)。このように順次的に自動微分を計算していく手法をフォワードモード微分と呼ぶ。

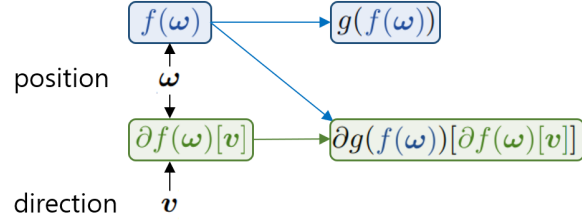


図3.1 $g \circ f$ の方向微分 $\partial g(f(\omega))[\partial f(\omega)[v]]$ 、JVP の計算グラフ。

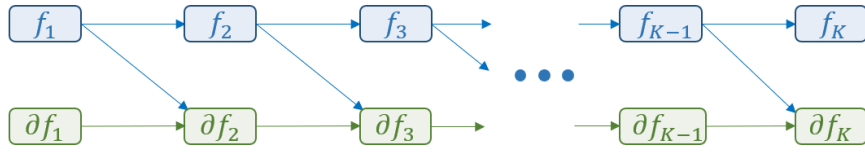


図3.2 フォワードモード微分。

もう少し具体的に $\omega \in \mathbb{R}^N, \mathbf{f} \in \mathbb{R}^M$ 、として、複数の関数の合成関数である $\mathbf{f}(\omega)$ を考える。JVP を用いて、Jacobian、

$$J_{\mathbf{f}, \omega} \equiv \frac{\partial \mathbf{f}}{\partial \omega} = \begin{pmatrix} \frac{\partial f_1}{\partial \omega_1} & \frac{\partial f_1}{\partial \omega_2} & \cdots & \frac{\partial f_1}{\partial \omega_N} \\ \frac{\partial f_2}{\partial \omega_1} & \frac{\partial f_2}{\partial \omega_2} & \cdots & \frac{\partial f_2}{\partial \omega_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_M}{\partial \omega_1} & \frac{\partial f_M}{\partial \omega_2} & \cdots & \frac{\partial f_M}{\partial \omega_N} \end{pmatrix} \in \mathbb{R}^{M \times N} \quad (3.15)$$

を計算したい。 \mathbf{f} についての JVP、

$$\text{JVP}(\mathbf{f}, \mathbf{v}) = J_{\mathbf{f}, \omega} \mathbf{v} = \left(\frac{\partial f_1}{\partial \omega} \mathbf{v}, \dots, \frac{\partial f_M}{\partial \omega} \mathbf{v} \right)^\top \quad (3.16)$$

が計算できるとすると、 $\frac{\partial f_i}{\partial \omega_j}$ は、

$$\text{JVP}(\mathbf{f}, \mathbf{e}_j) = \frac{\partial \mathbf{f}}{\partial \omega_j} \quad (3.17)$$

の i 成分を取り出すことで計算できる。ここに \mathbf{e}_j は j 番目の単位ベクトル (j 番目の要素だけが1で残りが0のベクトル) である。そこで \mathbf{f} の JVP を、各関数の JVP から合成することを考えよう。

まず $\mathbf{f}(\omega) = \mathbf{f}_2(\mathbf{f}_1(\omega))$ という場合の連鎖がどうなるか考えよう。 \mathbf{f} の微分は $\omega_1 = \mathbf{f}_1$ と書くと、

$$J_{\mathbf{f}, \omega} = \frac{\partial \mathbf{f}}{\partial \omega} = \left(\frac{\partial \mathbf{f}_2(\omega_1)}{\partial \omega_1} \right) \left(\frac{\partial \omega_1}{\partial \omega} \right) = \left(\frac{\partial \mathbf{f}_2(\omega_1)}{\partial \omega_1} \right) \left(\frac{\partial \mathbf{f}_1(\omega)}{\partial \omega} \right) \quad (3.18)$$

$$= J_{\mathbf{f}_2, \omega_1} J_{\mathbf{f}_1, \omega} \quad (3.19)$$

である。つまり一般に $\mathbf{f}(\omega) = \mathbf{f}_L(\mathbf{f}_{L-1}(\cdots \mathbf{f}_2(\mathbf{f}_1(\omega)))$ の場合、

$$J_{\mathbf{f}, \omega} = J_{\mathbf{f}_L, \omega_{L-1}} \cdots J_{\mathbf{f}_2, \omega_1} J_{\mathbf{f}_1, \omega} \quad (3.20)$$

ということである。

式 (3.20) の右辺に \mathbf{v} をかければ

$$J_{\mathbf{f}, \boldsymbol{\omega}} \mathbf{v} = J_{\mathbf{f}_L, \boldsymbol{\omega}_{L-1}} \cdots J_{\mathbf{f}_2, \boldsymbol{\omega}_1} J_{\mathbf{f}_1, \boldsymbol{\omega}} \mathbf{v} \quad (3.21)$$

$$= J_{\mathbf{f}_L, \boldsymbol{\omega}_{L-1}} (\cdots J_{\mathbf{f}_3, \boldsymbol{\omega}_1} (J_{\mathbf{f}_2, \boldsymbol{\omega}_1} (J_{\mathbf{f}_1, \boldsymbol{\omega}} \mathbf{v})) \cdots) \quad (3.22)$$

となるが、これは式 (3.30) のことである。右 (小さい L の順) から JVP を適用していけば $J_{\mathbf{f}, \boldsymbol{\omega}} \mathbf{v}$ ひいては $J_{\mathbf{f}, \boldsymbol{\omega}}$ の各成分が求まることがわかる。フォワードモードに必要なのは $\text{JVP}(\mathbf{f}_1, \mathbf{v}), \text{JVP}(\mathbf{f}_2, \mathbf{v}), \dots, \text{JVP}(\mathbf{f}_L, \mathbf{v})$ であることが確認された。

3.1.2 ヤコビアン の 計算量 と メモリ 使用量

与えられた方向 $\mathbf{v} \in \mathbb{R}^D$ に対して JVP が計算できるとして、ヤコビアン $\partial \mathbf{f} \in \mathbb{R}^{M \times D}$ を求めるためには式 (3.3) より、

$$\partial_i \mathbf{f}(\boldsymbol{\omega}) = \partial \mathbf{f}(\boldsymbol{\omega})[\mathbf{e}_i] \in \mathbb{R}^M \quad (3.23)$$

のように各単位ベクトル $\mathbf{e}_i \in \mathbb{R}^D$ を $i = 1, 2, \dots, D$ について計算すればよい。つまり一回のヤコビアンの評価には D 個の JVP を必要とする。関数の合成一回分を 1 レイヤーとすると $\mathbf{f}^k \in \mathbb{R}^{D_k}$ から $\mathbf{f}^{k-1} \in \mathbb{R}^{D_{k-1}}$ への計算は $D_k \times D_{k-1}$ の計算コストがかかるので、連鎖の各レイヤーで D の数が不変とする ($D_k = D_{k-1} = D$) と、連鎖のレイヤー数を K のとき KD^3 の計算コストがかかる。

3.2 Vector Jacobian Product

JVP、つまり方向微分 $\partial \mathbf{f}(\boldsymbol{\omega})[\mathbf{v}]$ と双対となる \mathbf{u} との内積 $\langle \cdot, \cdot \rangle \in \mathbb{R}$ を考え、その随伴を Vector Jacobian Product (VJP) と呼ぶ。すなわち

$$\langle \mathbf{u}, \partial \mathbf{f}(\boldsymbol{\omega})[\mathbf{v}] \rangle = \langle \partial \mathbf{f}(\boldsymbol{\omega})^*[\mathbf{u}], \mathbf{v} \rangle \quad (3.24)$$

となる $\partial \mathbf{f}(\boldsymbol{\omega})^*[\mathbf{u}]$ が VJP である。 $\boldsymbol{\omega} \in \mathbb{R}^N$ がベクトルである時、Jacobian を $J \equiv \partial \mathbf{f} \in \mathbb{R}^{M \times N}$ と略記することで、上記の関係は

$$\mathbf{u}^T (J\mathbf{v}) = (J^T \mathbf{u})^T \mathbf{v} \quad (3.25)$$

となることから、VJP は

$$\partial \mathbf{f}(\boldsymbol{\omega})^*[\mathbf{u}] = J^T \mathbf{u} \quad (3.26)$$

となる。

ベクトルの内積は $\mathbf{a}^T \mathbf{b} = \text{tr}(\mathbf{b}^T \mathbf{a})$ のように trace を使ってもあらわされるが、行列の場合の内積も同様に trace を使って定義する。これはフロベニウス内積となる。

$$\langle A, B \rangle = \text{tr}(B^T A) = \sum_{i,j} A_{ij} B_{ij} \quad (3.27)$$

この定義は行列 A, B をベクトル化したものにベクトル内積を適用したものとも同値であることに注意する。

ω が行列の場合

1. $f(W) = Wx$

$f(W) = Wx \in \mathbb{R}^N$ の VJP は $\partial f(W)[V] = Vx$ であったので $u \in \mathbb{R}^N$ として、

$$\langle u, Vx \rangle = u^T(Vx) = \text{tr}((Vx)^T u) = \text{tr}(ux^T V^T) = \text{tr}(xu^T V) = \langle ux^T, V \rangle \quad (3.28)$$

より VJP は $\partial f(W)^*[u] = ux^T$ となる。ただし $\text{tr}(x^T y) = y^T x$, $\text{tr}(A) = \text{tr}(A^T)$ と trace の巡回性を用いた。

2. $f(W) = WX$

JVP は $\partial f(W)[V] = VX$ であったので

$$\langle U, VX \rangle = \text{tr}((VX)^T U) = \text{tr}(UX^T V^T) = \text{tr}(XU^T V) = \langle UX^T, V \rangle \quad (3.29)$$

より VJP は $\partial f(W)^*[U] = UX^T$ となる。

3.2.1 VJP の連鎖則

JVP の時と同様に、gradient の連鎖則 (3.13) を用いて

$$\partial(g \circ f)(\omega)^*[u] = \partial(g \circ f)(\omega)^T u = (\partial g(f(\omega)) \partial f(\omega))^T u = \partial f(\omega)^T \partial g(f(\omega))^T u \quad (3.30)$$

$$= \partial f(\omega)^*[\partial g(f(\omega))^T u] \quad (3.31)$$

となることより、

$$\partial(g \circ f)(\omega)^*[u] = \partial f(\omega)^*[\partial g(f(\omega))^*[u]] \quad (3.32)$$

を得る。図3.3に VJP の連鎖則の計算グラフを示す。JVP の時とは異なり、 f から g へと関数と微分を順次的に計算することはできない。

図3.4に示すように、VJP ではまず関数そのものの合成を行う forward path の計算を行い、それぞれの関数値をメモリに保存しておき、backward path で微分を計算していくときに利用する。このため JVP より多くのメモリを必要とする。

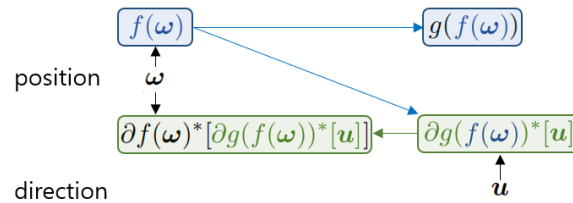


図3.3 $g \circ f$ の VJP の連鎖則の計算グラフ。

またもう少し具体的に ω がベクトルの場合の VJP を見てみよう。VJP は、

$$\text{VJP}(f, u) = (u^\top J)^\top = \left(u \cdot \frac{\partial f}{\partial \omega_1}, u \cdot \frac{\partial f}{\partial \omega_2}, \dots, u \cdot \frac{\partial f}{\partial \omega_N} \right)^\top \quad (3.33)$$

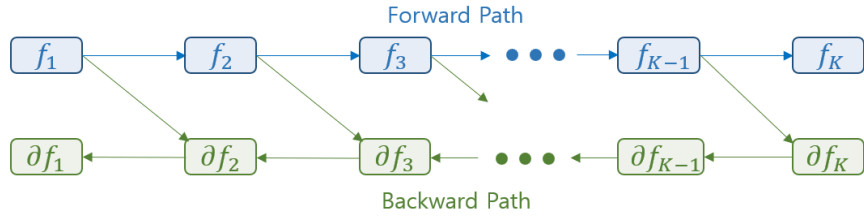


図3.4 リバースモード微分。

と定義される。この場合も同様に $\mathbf{u} = \mathbf{e}_i$ とおけば、

$$\text{VJP}(\mathbf{f}, \mathbf{e}_i) = (\mathbf{e}_i^\top J_{\mathbf{f}, \boldsymbol{\omega}})^\top = \left(\frac{\partial f_i}{\partial \omega_1}, \frac{\partial f_i}{\partial \omega_2}, \dots, \frac{\partial f_i}{\partial \omega_N} \right)^\top = \frac{\partial f_i}{\partial \boldsymbol{\omega}} = \nabla f_i \quad (3.34)$$

が得られるので、Jacobian 各成分が求まることがわかる。VJP を用いる連鎖則の計算法はリバースモードとよばれ、左から

$$\mathbf{u}^\top J_{\mathbf{f}, \boldsymbol{\omega}} = \mathbf{u}^\top J_{\mathbf{f}_L, \boldsymbol{\omega}_{L-1}} \cdots J_{\mathbf{f}_2, \boldsymbol{\omega}_1} J_{\mathbf{f}_1, \boldsymbol{\omega}} \quad (3.35)$$

$$= (\cdots ((\mathbf{u}^\top J_{\mathbf{f}_L, \boldsymbol{\omega}_{L-1}}) J_{\mathbf{f}_{L-1}, \boldsymbol{\omega}_1}) J_{\mathbf{f}_{L-2}, \boldsymbol{\omega}_1} \cdots) J_{\mathbf{f}_1, \boldsymbol{\omega}} \quad (3.36)$$

となり backward differentiation に必要なのは $\text{VJP}(\mathbf{f}_L, \mathbf{u})$, $\text{VJP}(\mathbf{f}_{L-1}, \mathbf{u})$, \dots , $\text{VJP}(\mathbf{f}_1, \mathbf{u})$ ということになる。

3.2.2 JVP と VJP の計算量・メモリ使用量の比較

3.2.3 小まとめ

	JVP (Jacobian-Vector Product)	VJP (Vector-Jacobian Product)
表記	$\partial f(\boldsymbol{\omega})[\mathbf{v}]$	$\partial f(\boldsymbol{\omega})^*[\mathbf{u}]$
連鎖則	$\partial(g \circ f)(\boldsymbol{\omega})[\mathbf{v}] = \partial g(f(\boldsymbol{\omega}))[\partial f(\boldsymbol{\omega})[\mathbf{v}]]$	$\partial(g \circ f)(\boldsymbol{\omega})^*[\mathbf{u}] = \partial f(\boldsymbol{\omega})^*[\partial g(f(\boldsymbol{\omega}))^*[\mathbf{u}]]$
$f(W) = W\mathbf{x}$	$\partial f(W)[V] = V\mathbf{x}$	$\partial f(W)^*[\mathbf{u}] = \mathbf{u}\mathbf{x}^\top$
$f(W) = WX$	$\partial f(W)[V] = VX$	$\partial f(W)^*[U] = UX^\top$

表3.1 JVP/VJP まとめ

3.3 JAX の自動微分をカスタマイズする

というわけで JAX に戻って、JVP, VJP を自分で定義してみよう。フォワードモードのカスタム自動微分は、JVP を定義する `jax.custom_jvp` を使用すれば良い。さて以下の例では

$$g(\mathbf{x}) = x_2 \sin x_1 \quad (3.37)$$

の JVP を実装する。みやすさのため $x = x_1$, $y = x_2$ とおいている。

```

1 from jax import custom_vjp
2 @custom_vjp
3 def g(x, y):
4     return jnp.sin(x) * y
5
6 @g.defjvp
7 def g_jvp(primals, tangents):
8     x, y = primals
9     ux, uy = tangents
10    dgdx=y * jnp.cos(x)
11    dgdy=jnp.sin(x)
12    primal_out = g(x, y)
13    tangent_out = dgdx * ux + dgdy * uy
14    return primal_out, tangent_out

```

図3.5に JVP とインプット、アウトプットの対応を示す。tangent_out = dg1dx * ux + dg1dy * uy とおいているが、これが JVP の微分更新部分になっている。つまり、

$$\text{JVP}(g, \mathbf{u}) = J_{g, \mathbf{x}} \mathbf{u} = \frac{\partial g}{\partial \mathbf{x}} \mathbf{u} = \frac{\partial g(x, y)}{\partial x} u_1 + \frac{\partial g(x, y)}{\partial y} u_2 \quad (3.38)$$

$$= (y \cos x) u_x + (\sin x) u_y \quad (3.39)$$

となっている。

Jacobian-Vector Product (JVP)

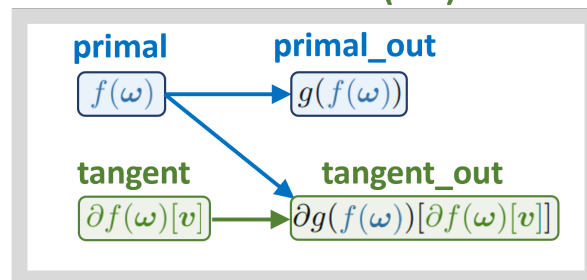


図3.5 JVP の実装。

リバースモードのカスタム自動微分は、VJP を定義する `jax.custom_vjp` を使用すれば良い。図3.4のようにリバースモードでは、Forward path (fwd) の部分と Backward path (bwd) の部分の二つのメソッドを定義する必要がある。また forward path では、後に backward path で微分を計算するために必要な情報を出力しておく。これを residuals と呼ぶ。この部分は backward path で微分を計算するのに必要ならば、どのような数学的情報でも良いが、メモリ使用量に直接影響するので、最小限にしたほうがよい。

以下では fwd と bwd の実装例を示す。また図3.6に VJP の実装部分のインプット・アウトプットの対応を示す。fwd および bwd は最後に `defvjp` により VJP として統合して定義される。ちなみにこの例では y と

`cos_x` を別々に保存しているが、これは掛けた状態で保存したほうがメモリ使用量的には効率が良い。

```

1 from jax import custom_vjp
2
3 @custom_vjp
4 def g(x, y):
5     return jnp.sin(x) * y
6
7 def g_bwd(residuals, u):
8     cos_x, sin_x, y = residuals
9     return (y * cos_x * u, sin_x * u)
10
11 def g_fwd(x, y):
12     residuals = (jnp.cos(x), jnp.sin(x), y)
13     primal_out = g(x, y)
14     return primal_out, residuals
15
16 g.defvjp(g_fwd, g_bwd)

```

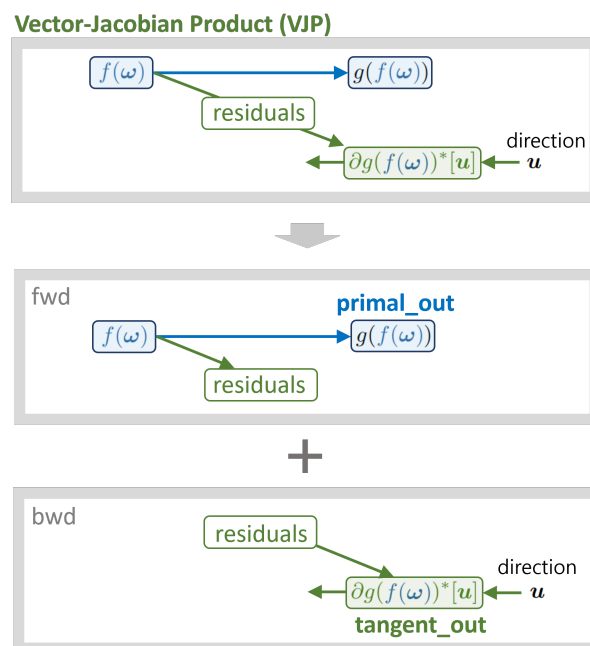


図3.6 VJP の実装。

3.4 Newton-Raphson 法で実装された陰関数の VJP を定義する

一般のプログラム中には陽的な関数ばかりでなく、反復的に解を求める陰関数も使われる。このような場合の微分は、自分で JVP なり VJP を定義する必要がある。ここでは例題として、力学の二体問題中の問題より、eccentricity e と Mean anomaly M から eccentric anomaly E を求める問題を考えよう。とはいえ、力学に興味が無くても問題設定は理解できる。eccentric anomaly は

$$f(E) = E - e \sin E - M = 0 \quad (3.40)$$

を解くことで E が求まる、というだけである。この非線形方程式 $f(E) = 0$ は数値的に求めることが必要である。このための方法の一つとして Newton-Raphson 法がある。Newton-Raphson 法は図3.7にしめすように、まず初期値 E_1 から初めて、 $f(E_1)$ に接する直線を求め、その直線と原点との交点を解析的に求め E_2 とする。この手続きを n 回繰り返すという手順である。 i 番目の手続きでは接線は

$$y = f'(E_i)(E - E_i) + f(E_i) \quad (3.41)$$

となることから

$$E_{i+1} = E_i - \frac{f(E_i)}{f'(E_i)} \quad (3.42)$$

となる。この手続きを収束条件の判定値を ϵ として $|E_{i+1} - E_i| < \epsilon$ となるまで繰り返し、条件を満たした E_i を近似解とすればよい。

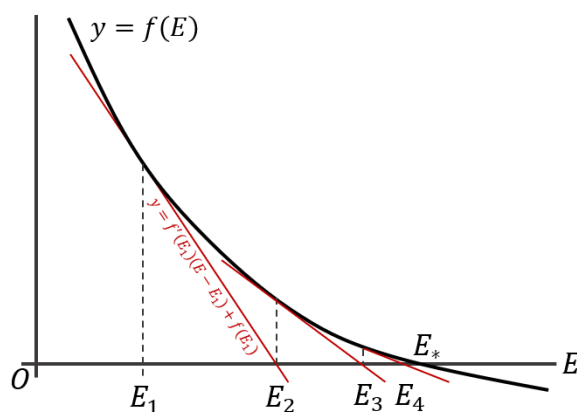


図3.7 Newton-Raphson 法。

さて、いま E は e と M の関数とみなすことができる。すなわち

$$E = E(e, M) \quad (3.43)$$

を上記の Newton-Raphson 法で求める関数を JAX で実装するには、例えば、以下のようにすればよい。ここで e , M 以外にも、Newton-Raphson 法の初期値 E_{ini} も引数に入っている。

```
1 from jax.lax import while_loop
2 import jax.numpy as jnp
```

```

3 from jax import grad
4
5 def f(E, e, M):
6     return E - e*jnp.sin(E) - M
7
8 dfdE = grad(f, argnums=0)
9
10 @custom_vjp
11 def eccentric_anomaly_newton_raphson(e, M, Eini):
12     def cond_fun(carry):
13         E_prev, E = carry
14         return jnp.abs(E - E_prev) > 1e-6
15
16     def body_fun(carry):
17         _, E = carry
18         E_new = E - f(E, e, M) / dfdE(E, e, M)
19         return E, E_new
20
21     _, E_star = while_loop(cond_fun, body_fun, (jnp.inf, Eini))
22     return E_star

```

さて、この関数は while_loop が入っているので reverse-mode で微分が通らない。そこで VJP を自分で定義してやる必要がある。微分は自分で用意する必要があるが、陰関数定理より

$$\frac{\partial E}{\partial e} = -\frac{\partial f / \partial e}{\partial f / \partial E} = \frac{\sin E}{1 - e \cos E} \quad (3.44)$$

$$\frac{\partial E}{\partial M} = -\frac{\partial f / \partial M}{\partial f / \partial E} = \frac{1}{1 - e \cos E} \quad (3.45)$$

である。これより forward path (fwd) では、 $\sin E$ と $1 - e \cos E$ を保存しておいて、backward path (bwd) で用いればよいことが分かる。つまり全体として以下のようにすればよい。

```

1 import jax.numpy as jnp
2 from jax import grad
3
4 def f(E, e, M):
5     return E - e*jnp.sin(E) - M
6
7 dfdE = grad(f, argnums=0)
8 dfde = grad(f, argnums=1)
9 dfdM = grad(f, argnums=2)
10

```

```

11 from functools import partial
12 from jax import custom_vjp
13 from jax.lax import while_loop
14
15 @custom_vjp
16 def eccentric_anomaly_newton_raphson(e, M, Eini):
17     def cond_fun(carry):
18         E_prev, E = carry
19         return jnp.abs(E - E_prev) > 1e-6
20
21     def body_fun(carry):
22         _, E = carry
23         E_new = E - f(E, e, M) / dfdE(E, e, M)
24         return E, E_new
25
26     _, E_star = while_loop(cond_fun, body_fun, (jnp.inf, Eini))
27     return E_star
28
29
30 def eccentric_anomaly_newton_raphson_fwd(e, M, Eini):
31     E_star = eccentric_anomaly_newton_raphson(e, M, Eini)
32     return E_star, (jnp.sin(E_star), 1.0 - e * jnp.cos(E_star))
33
34
35 def eccentric_anomaly_newton_raphson_bwd(residuals, u):
36     sin_E_star, one_minus_cos_E_star = residuals
37     return (sin_E_star * u / one_minus_cos_E_star, u / one_minus_cos_E_star, 0.0)
38
39
40 eccentric_anomaly_newton_raphson.defvjp(eccentric_anomaly_newton_raphson_fwd,
41                                         eccentric_anomaly_newton_raphson_bwd)

```

$\partial E / \partial M$ が計算できるか確認してみよう。

```

1 from jax import jit, vmap
2
3 dE_de = grad(eccentric_anomaly_newton_raphson, argnums=0)
4 dE_dM = grad(eccentric_anomaly_newton_raphson, argnums=1)
5 vmap_eccentric_anomaly = jit(vmap(eccentric_anomaly_newton_raphson, in_axes=(None, 0, 0)))
6 vmap_dE_de = jit(vmap(dE_de, in_axes=(None, 0, 0)))

```



```

7 | vmap_dE_dM = jit(vmap(dE_dM, in_axes=(None, 0, 0)))
8 |
9 | etest = 0.9
10 | M_values = jnp.linspace(0, 2*jnp.pi-0.001, 100)
11 | E_stars = vmap_eccentric_anomaly(etest, M_values, M_values)
12 | dE_stardM = vmap_dE_dM(etest, M_values, M_values)

```

プロットしてみると、ちゃんと微分が計算できている。

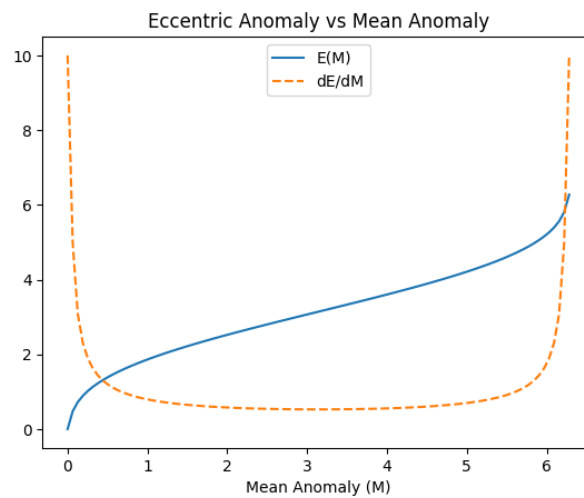


図3.8 $E(M)$ および $\partial E/\partial M$ 。

第 4 章

JAX プログラミングの様々なテクニック

4.1 XLA の再コンパイルを防ぐ

JAX/XLA は、実行時にコンパイルを行う。このコンパイルは通常、最初の一回のみであるが、関数のハッシュ値が変わると再コンパイルが行われてしまう。これは

- 同じ関数を繰り返し使用する
- 関数実行時間に対してコンパイルの時間が同程度か長い

という二つの条件が満たされると、再コンパイルの時間が実行時間を律速しうる（図4.1）。このような状況が起きているかどうかは、`prefetto` のようなメモリプロファイラーを用いれば分かる。対策としては `jit` や `scan`, `vmap` の位置を変更するなどがあるが現状手探りで探すしかない。

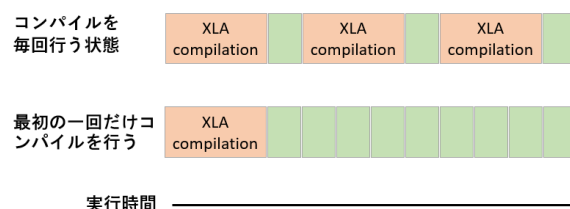


図4.1 コンパイルの実行時間。

4.2 reverse-mode 微分と checkpoint

checkpoint を置くことで、その時点での reverse-mode の前進方向のメモリを開放することができる。ややこしいのでまた後で書く。

4.3 `jax.numpy.array` の範囲外アサイン

python の `numpy.array` は未定義のインデックスにアクセスしようとするとき `IndexError` が出される。しかし `jax.numpy.array` ではこのようなことは起きず、なんらかの値が返されてしまう。直接、範囲外のインデックスを入れると、通常は最終の要素を返すようだが、これも厳密には未定義動作である。一方、`jnp.take` を用

いと、まったく関係ない値が返される。つまり、JAX ではあらかじめインデックスが有効範囲内にあることをチェックすることが必要である。

例えば

```
1 import jax.numpy as jnp
2 a = jnp.array([1,2,3])
3 a[0],a[1],a[2],a[3],a[4]
```

では

```
1 (Array(1, dtype=int32),
2  Array(2, dtype=int32),
3  Array(3, dtype=int32),
4  Array(3, dtype=int32),
5  Array(3, dtype=int32))
```

が返ってくる（おそらく保証はされていない）。一方、

```
1 jnp.take(a, 3)
```

では、

```
1 Array(-2147483648, dtype=int32)
```

のようなとんでもない値が返ってくる！

4.4 jax.numpy.array のインデックス操作

jax.numpy は直接のインデックスアサインはできない。以下のコードでは5成分のゼロアレイを用意して、インデックス 1,2 にそれぞれ 10,3 を代入しようとしている。しかし、最後の行はエラーになる。

```
1 a=jnp.zeros(5)
2 indx=jnp.array([1,2])
3 val=jnp.array([10.0,3.0])
4 a[indx]=val #エラー
```

代わりに jax.numpy.at.set を用いることができる。

```
1 a=a.at[indx].set(val)
2 # DeviceArray([ 0., 10.,  3.,  0.,  0.], dtype=float32)
```

とすればよい。

ここですでにある配列のインデックスアサインをして足すには at.add を用いる。

```

1 a.at[jnp.array([0,1])].add(jnp.array([1.,2.]))
2 # DeviceArray([ 1., 12.,  3.,  0.,  0.], dtype=float32)

```

注意：jax.opt.index などの関数は廃止されました。

4.5 scan 中のインデクシング

dynamic_slice は jax.array の一部分をスライスするのに用いられる。

```

1 from jax.lax import scan
2 from jax.lax import dynamic_slice
3 import jax.numpy as jnp
4
5 N = 30
6 M = 10
7 a = jnp.arange(N)
8
9 print(dynamic_slice(a, (10,), (M,)))

```

は

```

1 [10 11 12 13 14 15 16 17 18 19]

```

を返してくる。

この dynamic_slice を用いて、scan の中で slice することができる。ここでは配列 a を順にスライスしたものに c[i] をかけて、append したものを得る例をあげる

```

1 c = jnp.array([1,3,7])
2
3 def scan_slice():
4
5     def f(icarry,factor):
6         element = dynamic_slice(a, (icarry,), (M,))
7         return icarry+M, element*factor
8
9     _, result = scan(f, 0, c)
10
11     return result

```

この関数を呼ぶと

```
1  [[ 0  1  2  3  4  5  6  7  8  9]
2   [ 30 33 36 39 42 45 48 51 54 57]
3   [140 147 154 161 168 175 182 189 196 203]]
```

を返す。

第 5 章

JAX で最適化

ちょっと古いです。JAXOpt は冷温停止状態で、Optax に統合されるようです。今後更新予定。

5.1 jax.experimental

最適化は自動微分の代表的な応用先である。多くの最適化手法はもとを正せば Gradient Descent

$$\mathbf{r}^{(k)} = \mathbf{r}^{(k-1)} - \gamma \frac{\partial}{\partial \mathbf{r}} f(\mathbf{r}) \quad (5.1)$$

に似た方法で更新していくからであり、つまり目的関数の微分 $f'(\mathbf{r})$ が必要だからである。つまり自動微分が通るように構成しておけば、自分で微分を計算する必要はない。

5.1.1 JAX の optimizer で最適化を楽しむ

JAX には `jax.experimental.optimizer` という最適化モジュールがある。有名どころの optimizer はここのもを使って手軽に最適化できる。例として Booth function

$$f(x, y) = (x + 2y - 7)^2 + (2x + y - 5)^2 \quad (5.2)$$

を目的関数として、これを最小化する (x, y) を求めてみよう。

```
1 #Booth function
2 def booth(x,y):
3     f=(x+2.*y-7. )**2 + (2.*x+y-5. )**2
4     return f
```

```
1 # 目的関数
2 def objective(r):
3     f=booth(r[0],r[1])
4     return f
```

`jax.experimetnal.optimizer` は `opt_init` (初期化), `opt_update` (アップデート), `get_params` (パラメタをゲットする) の 3 セットを用いて最適化を行う。ステップを以下のように定義する。

```

1 from jax import value_and_grad
2 def step(t, opt_state):
3     value, grads = value_and_grad(objective)(get_params(opt_state))
4     opt_state = opt_update(t, grads, opt_state)
5     return value, opt_state

```

`value_and_grad(f)(p)` は引数 $x = p$ での関数 $f(x)$ の値と微分値 $f'(x)$ を返す。t は step index（整数値）である。微分値と step index、前回の状態 (`opt_state`) を `opt_update` に渡すことで、次のステップにアップデートする。初期値を `r0=(x0,y0)` として、Nstep 回ステップを回して、途中結果を保存するには

```

1 def doopt(r0,opt_init,get_params,Nstep):
2     opt_state = opt_init(r0)
3     xtraj=[r0[0]]
4     ytraj=[r0[1]]
5
6     for t in range(Nstep):
7         value, opt_state = step(t, opt_state)
8         p=get_params(opt_state)
9         xtraj.append(p[0])
10        ytraj.append(p[1])
11    return xtraj, ytraj, p

```

のようにすればよい。optimizer としては ADAM と SGD を試す。

```

1 from jax.experimental import optimizers
2 #ADAM
3 opt_init, opt_update, get_params = optimizers.adam(1e0)
4 r0 = jnp.array([7.5,-7.5]) # 初期値
5 xtraj1, ytraj1, p=doopt(r0,opt_init,opt_update,get_params,100)

```

もしくは

```

1 #SGD
2 opt_init, opt_update, get_params = optimizers.sgd(1e-1)
3 r0 = jnp.array([7.5,-7.5])
4 xtraj2, ytraj2, p=doopt(r0,opt_init,get_params,100)

```

のような感じである。実行結果を表示しよう。

```

1 import matplotlib.pyplot as plt
2 import jax.numpy as jnp

```

```

3 import seaborn as sns
4 plt.style.use('bmh')
5
6 #grid for contour
7 Nx=1000;Ny=1000
8 x=jnp.linspace(-5,15,Nx)
9 y=jnp.linspace(-10,10,Ny)
10 g=booth(x[:,None]*jnp.ones((Nx,Ny)),y[None,:]*jnp.ones((Nx,Ny)))
11
12 fig=plt.figure(figsize=(6,6))
13 plt.plot(xtraj1,ytraj1,lw=0.75,c="C1",label="ADAM")
14 plt.plot(xtraj2,ytraj2,lw=0.75,c="C0",ls="dashed",label="SGD")
15 plt.contour(x,y,jnp.log(g.T),levels=100,cmap="gray",alpha=0.3)
16 plt.xlabel("x")
17 plt.ylabel("y")
18 plt.legend()

```

結果は、図5.1のようになる。これは ADAM が SGD より悪いというわけではなく step_size を大きめ

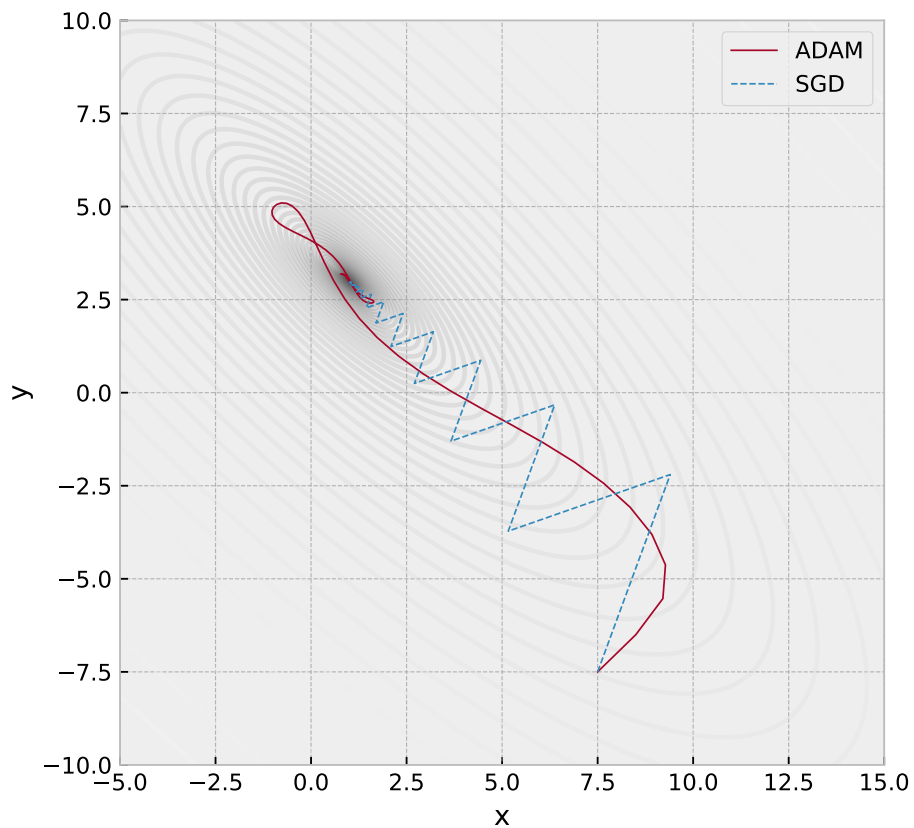


図5.1

に設定したので慣性が効いているのが見えているというわけである。Wikipedia に[最適化のテスト関数](https://en.wikipedia.org/wiki/Test_functions_for_optimization) (https://en.wikipedia.org/wiki/Test_functions_for_optimization)が多数挙げられているので試してみよう。全然、大域最小値に収束しないこと請け合いである。

5.2 JAXopt

[JAXopt](#)は JAX を用いた最適化のためのパッケージである。特徴として目を引くのは

- Differentiable: optimization problem solutions can be differentiated with respect to their inputs either implicitly or via autodiff of unrolled algorithm iterations.

すなわち、解自体を最適化問題のパラメタで微分できるというところにある。すばらしい。

以下では $f(x; a) = (x - \sin(a))^2$ を x について最適化してみよう。もちろん答えは $x_* = \sin(a) \equiv g(a)$ となるはずである。

```
1 import jaxopt
2 def f(x,a):
3     return (x-jnp.sin(a))**2
4 gd = jaxopt.GradientDescent(fun=f, maxiter=500)
```

これだけである。解 $x_* = g(a)$ を以下のように定義すればよい。

```
1 def g(a):
2     res = gd.run(init_params=np.random.normal(0.0,1.0), a=a)
3     params, state = res
4     return params
```

これで、

```
1 g(0.1) #-> DeviceArray(0.09983341, dtype=float32)
```

となり解が出る。さらに $g(a)$ を a で微分することができる。

```
1 from jax import grad
2 dg=grad(g)
3 dg(-0.2) #-> DeviceArray(0.9800666, dtype=float32)
```

これは $\frac{\partial}{\partial a} x_{\min} = \frac{\partial}{\partial a} g(a) = \cos(a)$ のことであり、たしかにそうになっている。

第 6 章

pytree/tree-math (TBD)

まだ書いてない

tree_map

第 7 章

Numpyro でマルコフ鎖モンテカルロ・シミュレーションをする

NumPyro の MCMC は HMC-NUTS(Hamiltonian Monte Carlo – No-U-Turn Sampling) が基本である。HMC は勾配情報を用いた高受容率の MCMC であり、フィットするモデルのパラメタ数 D が多い場合にはランダムメトロポリス・ヘイスティング法より効率がよい。理論的には収束までの時間が HMC では $D^{5/4}$ 、ランダム MH では D^2 に比例するとされている [5]。ここでは HMC の説明はしないが、日本語では「ゼロからできる MCMC」(KF 理工学専門書) が筆者のおすすめである。HMC では計算に Leap-frog を使い、モデルのパラメタによる微分が必要である。numpyro [6] を用いて HMC-NUTS をしよう。numpyro は JAX を用いた Probabilistic Programming Language である。つまり JAX の自動微分が通るようにモデルを定義しないといけない。

7.1 曲線フィット

正弦曲線から生成されるデータにガウスノイズが加わったものを考えよう。すなわち

$$y = \sin(x + \phi) + \epsilon \quad (7.1)$$

$$\epsilon \sim \mathcal{N}(0, \sigma) \quad (7.2)$$

のような場合である。ここに $\epsilon \sim p$ というのは確率変数 ϵ が確率密度 p から生成されるということを意味している。また $\mathcal{N}(\mu, \sigma)$ は平均 μ 、標準偏差 σ の正規分布のことである。さてデータから位相 ϕ とノイズレベル σ を推定したい。まず模擬データを以下のように生成しておく。

```
1 import numpy as np
2 np.random.seed(32)
3 phase=0.5
4 sign=0.3
5 N=20
6 x=np.sort(np.random.rand(N))*4*np.pi
7 y=np.sin(x+phase)+np.random.normal(0,sign,size=N)
```

このデータ (x,y) は図7.1のような感じとなっている。

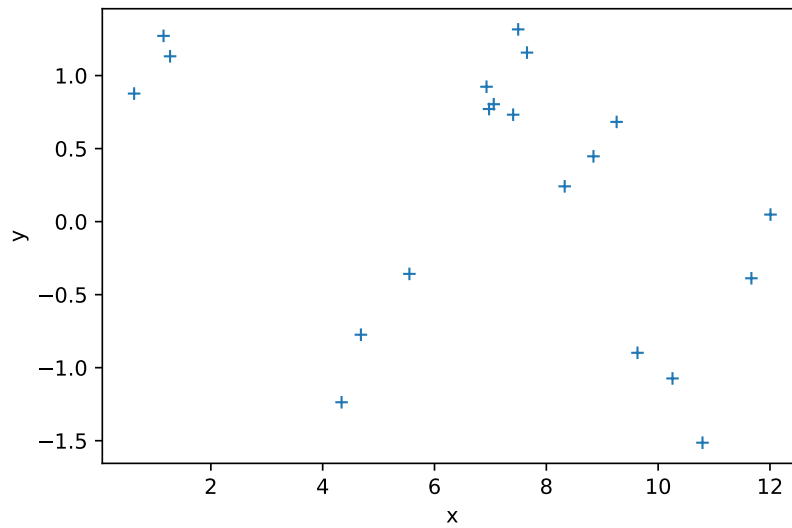


図7.1

さてこのデータに対し確率モデルを設定し HMC-NUTS で事後分布を求めよう。確率モデルとしては式 (7.1,7.2) また事前分布として

$$\phi \sim \mathcal{U}[-\pi, \pi] \quad (7.3)$$

$$\sigma \sim \text{Exponential}(1) \quad (7.4)$$

としたものを用いる。この確率モデルを numpyro を使って以下のように定義する。

```

1 import jax.numpy as jnp
2 import numpyro
3 import numpyro.distributions as dist
4
5 def model(x,y):
6     phase = numpyro.sample('phase', dist.Uniform(-1.0*jnp.pi, 1.0*jnp.pi))
7     sigma = numpyro.sample('sigma', dist.Exponential(1.))
8     mu=jnp.sin(x+phase)
9     numpyro.sample('y', dist.Normal(mu, sigma), obs=y)
```

この例からわかるように、確率分布は dist で指定されたものを、numpyro の sample 内で定義することにより確率変数として扱われる。最後のデータ y の部分は式 (7.1,7.2) が

$$y \sim \mathcal{N}(\sin(x + \phi), \sigma) \quad (7.5)$$

と書き直されることよりわかる。さてこれを numpyro の HMC-NUTS にかける。

```

1 from jax import random
```

```

2 from numpyro.infer import MCMC, NUTS
3 rng_key = random.PRNGKey(0)
4 rng_key, rng_key_ = random.split(rng_key)
5 num_warmup, num_samples = 1000, 2000

```

HMC-NUTS は最初にウォームアップで調整する。ウォームアップのチェーンの数と本番のチェーンの数がそれぞれ `num_warmup`, `num_samples = 1000, 2000` で指定されている。さて MCMC を実行する。

```

1 kernel = NUTS(model)
2 mcmc = MCMC(kernel, num_warmup=num_warmup, num_samples=num_samples)
3 mcmc.run(rng_key_, x=x, y=y)
4 mcmc.print_summary()

```

最後の `print_summary()` で以下のような事後分布サンプリングの要約が出力されたと思う。成功だ。

```

1 mean std median 5.0% 95.0% n_eff r_hat
2 phase 0.34 0.12 0.34 0.15 0.56 1363.51 1.00
3 sigma 0.34 0.06 0.33 0.24 0.42 1206.21 1.00
4 Number of divergences: 0

```

以下では結果のチェックに `arviz` を用いよう。MCMC は結果の健全性のチェックが重要である。`arviz` は標準的なチェックが兼ね備えられているので便利である。

```

1 import arviz
2 arviz.plot_trace(mcmc, var_names=["phase", "sigma"])

```

上では図7.2のように、パラメタの周辺分布（左）とチェーンが表示される。

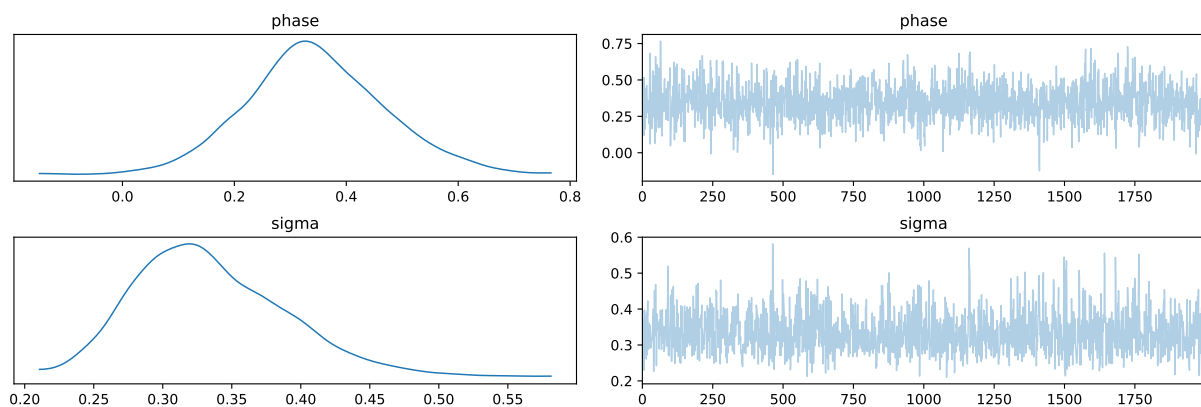


図7.2

いわゆるコーナープロットで事後分布の同時確率密度をチェックしよう^{*1}。

*1 2020/12月現在 `arviz` の develop 版を用いている。

```

1 refs={};refs["sigma"]=sign;refs["phase"]=phase
2 arviz.plot_pair(arviz.from_numpyro(mcmc),kind='kde',
3     divergences=False,marginals=True,reference_values=refs,
4     reference_values_kwargs={'color':"red", "marker":"o", "markersize":12})

```

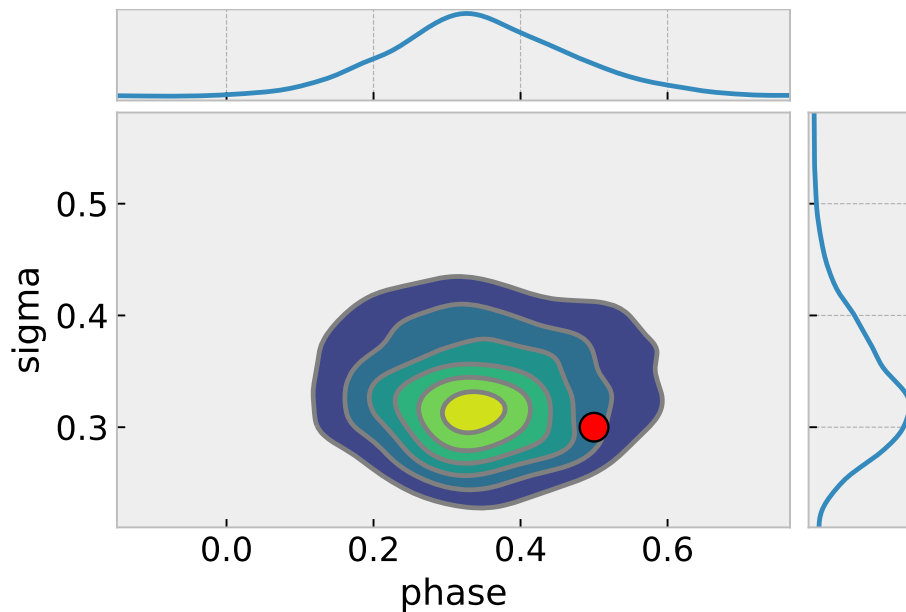


図7.3 赤い点は refs で指定されるインプットの値。

位相 ϕ および σ のサンプリングは

```

1 posterior_phase = mcmc.get_samples()['phase']
2 posterior_sigma = mcmc.get_samples()['sigma']

```

のようにして得ることができる。この事後分布サンプリングを用いて、予測分布のサンプリングをつくろう。
numpyro.infer の Predictive を用いることで簡単に予測分布をサンプルできる。

```

1 from numpyro.infer import Predictive
2 pred = Predictive(model,{'phase':posterior_phase,'sigma':posterior_sigma},
3     return_sites=["y"])
4 x_ = jnp.linspace(0,4*jnp.pi,1000)
5 predictions = pred(rng_key_,x=x_,y=None)

```

信頼分布として Highest posterior density interval(HPDI) を 90% 区間で計算しよう。

```

1 from numpyro.diagnostics import hpdi
2 mean_muy = jnp.mean(predictions["y"], axis=0)

```

```
3 hpdi_muy = hpdi(predictions["y"], 0.9)
```

最後にデータとともにプロットしよう。もちろんプロットの見た目を良くするだけの seaborn は使わなくても良い。結果は図7.4である。

```
1 import seaborn as sns
2 plt.style.use('bmh')
3
4 fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(6, 3))
5 ax.plot(x,y,"+",color="black")
6 ax.plot(x_,mean_muy,color="C0")
7 ax.fill_between(x_, hpdi_muy[0], hpdi_muy[1], alpha=0.3, interpolate=True,color="C0")
8 plt.xlabel("x")
9 plt.ylabel("y")
```

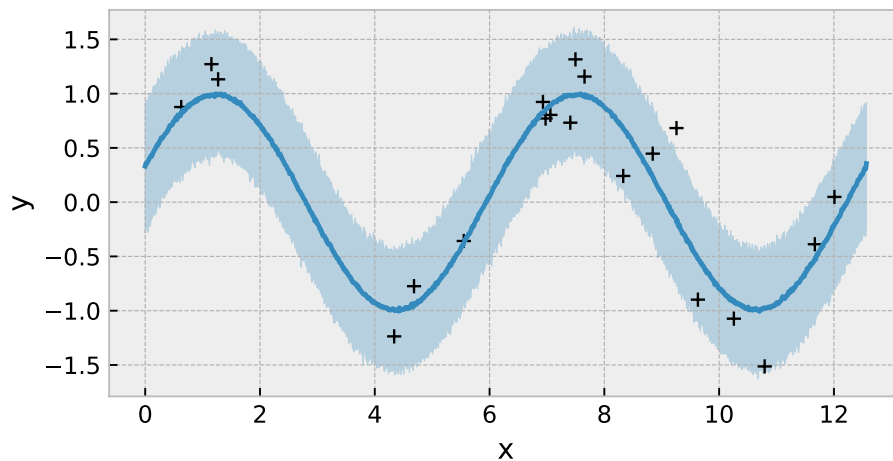


図7.4。

モデルに対する Prior のサンプリングを確認する

handlers を用いるとモデルからサンプリングできる。ただし、インプットに y は必要ないことに注意。

```
1 def modelc(x):
2     phase = numpyro.sample('phase', dist.Uniform(-1.0*jnp.pi, 1.0*jnp.pi))
3     sigma = numpyro.sample('sigma', dist.Exponential(1.))
4     mu=jnp.sin(x+phase)
5     d=numpyro.sample('y', dist.Normal(mu, sigma))
```

```
6     return d
```

これを `handlers.seed` を用いてサンプリングできる。

```
1 from numpyro import handlers
2 d = handlers.seed(modelc, rng_seed=9)(x)
```

対数尤度からサンプリングする

上の例では `model(x,y)` の中で `numpyro.sample` を用いて、データの確率分布を直接指定した。しかし場合によっては対数尤度からサンプリングしたいこともある。この場合、

```
1 def model(x,y):
2     phase = numpyro.sample('phase', dist.Uniform(-1.0*jnp.pi, 1.0*jnp.pi))
3     sigma = numpyro.sample('sigma', dist.Exponential(1.))
4     mu=jnp.sin(x+phase)
5     loglikelihood=-0.5*(y-mu)**2/sigma**2-jnp.log(sigma)
6     numpyro.factor("loglike", loglikelihood)
```

のようにして指定することもできる。

7.1.1 複数の種類のデータがある場合

共通のパラメタを変数にもつ2つの種類のデータがある場合を考えよう。例として

$$y_1 = \sin(x + \phi) + \epsilon \quad (7.6)$$

$$y_2 = \cos(x + \phi) + \epsilon \quad (7.7)$$

$$\epsilon \sim \mathcal{N}(0, \sigma) \quad (7.8)$$

という場合を考える。例えば図7.5のようなデータを解析する。

モデルは前節と同様に

$$y_1 \sim \mathcal{N}(\sin(x + \phi), \sigma) \quad (7.9)$$

$$y_2 \sim \mathcal{N}(\cos(x + \phi), \sigma) \quad (7.10)$$

と考えることができるので、`numpyro` のモデルは以下のように書ける。

```
1 def model(x1,x2,y1,y2):
2     phase = numpyro.sample('phase', dist.Uniform(-1.0*jnp.pi, 1.0*jnp.pi))
3     sigma = numpyro.sample('sigma', dist.Exponential(1.))
4     mu1=jnp.sin(x1+phase)
5     mu2=jnp.cos(x2+phase)
6     numpyro.sample('y1', dist.Normal(mu1, sigma), obs=y1)
7     numpyro.sample('y2', dist.Normal(mu2, sigma), obs=y2)
```

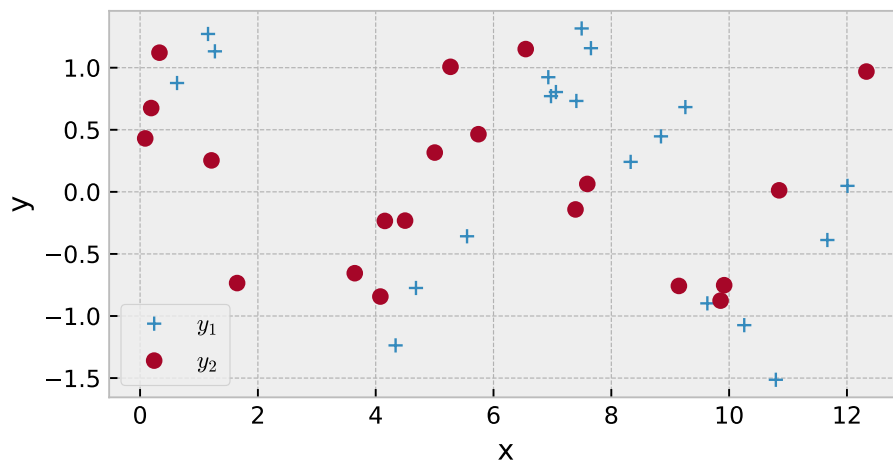



図7.5。

あとはほとんど前節と同じだが、MCMC 実行時に

```
1 mcmc.run(rng_key_, x1=x1, x2=x2, y1=y1, y2=y2)
```

のように変数の数を合わせるのと、予測分布を描くときにも以下のように2変数にする注意が必要である。

```
1 pred = Predictive(model,{'phase':posterior_phase,'sigma':posterior_sigma}
2     ,return_sites=["y1","y2"])
3 x1_ = jnp.linspace(0,4*jnp.pi,1000)
4 x2_ = jnp.linspace(0,4*jnp.pi,1000)
5 predictions = pred(rng_key_,x1=x1_,x2=x2_,y1=None,y2=None)
```

予測からもそれぞれ平均・HPDIを計算することにより、

```
1 mean_muy1 = jnp.mean(predictions["y1"], axis=0)
2 hpdi_muy1 = hpdi(predictions["y1"], 0.9)
3 mean_muy2 = jnp.mean(predictions["y2"], axis=0)
4 hpdi_muy2 = hpdi(predictions["y2"], 0.9)
```

最終的に、以下のようなプロットで

```
1 fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(6, 3))
2 ax.plot(x1,y1,"+",color="C0")
3 ax.plot(x1_,mean_muy1,color="C0")
4 ax.fill_between(x1_, hpdi_muy1[0], hpdi_muy1[1], alpha=0.3,
5     interpolate=True,color="C0")
6 ax.plot(x2,y2,"o",color="C1")
```

```

7 ax.plot(x2_,mean_muy2,color="C1")
8 ax.fill_between(x2_, hpdi_muy2[0], hpdi_muy2[1], alpha=0.3,
9                 interpolate=True,color="C1")
10 plt.xlabel("x")
11 plt.ylabel("y")

```

図7.6のように同時フィットすることができる。

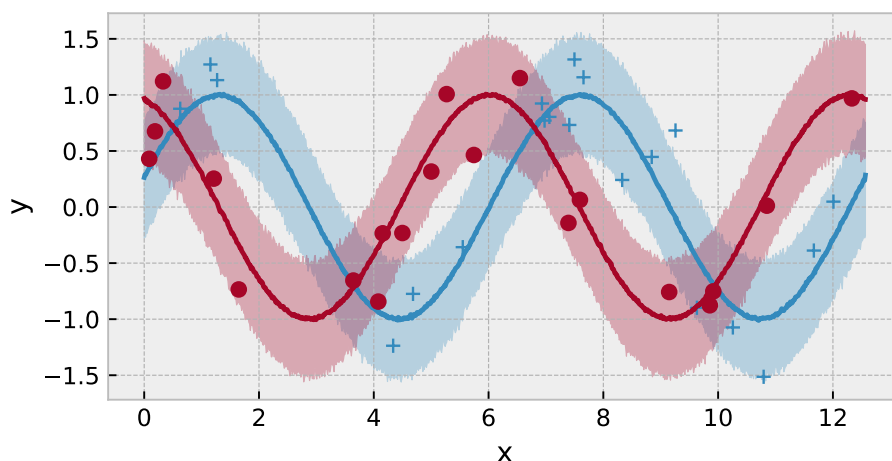


図7.6。

7.1.2 モデル内の関数がグリッドデータで与えられている時

世の中はモデルは往々にして美しい数学関数だけで記述されるわけではない。複雑なモデルになるほど、実験データからもとまる関係などが入ってくる。このような場合でも補間関数を用いることによって対処できる。一例として、図7.7の上パネルのようなデータ点（黒点）からモデルを構成することを考えよう。

このデータ点 x_a , y_a を `jax.numpy.interp` で補完した関数を以下のように定義する。

```

1 def f(x):
2     return jnp.interp(x,x_a,y_a)

```

すると、 $f(x)$ は図7.7の上パネル実線のように、データのある部分は良く補間されているのがわかる。さてこの $f(x)$ に対しても自動微分は通るので

```

1 from jax import grad,vmap
2 g=vmap(grad(f))

```

のように定義すれば、微分を計算できる。図7.7のしたパネルはそれを表示したもので、具体的には

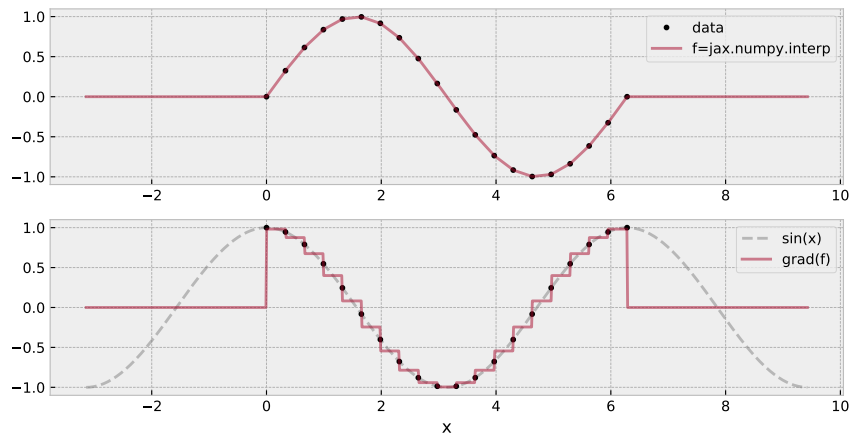


図7.7。

```

1 xx=np.linspace(-np.pi,3*np.pi,1000)
2 plt.plot(xx,np.cos(xx),ls="dashed",label="sin(x)",color="gray",alpha=0.5)
3 plt.plot(x,np.cos(x),".",color="black")
4 plt.plot(xx,g(xx),label="grad(f)",color="C1",alpha=0.5)
5 plt.legend()
6 plt.xlabel("x")

```

で、作成されている。そう、実は（というほどでもないが）このデータ点は $\sin(x)$ から作成されているのだが、そのことは知らないとして話を進めよう。補間関数の自動微分は、データ点の離散化に対応してカクカクしているものの、おおむね真のモデル $\partial_x \sin x = \cos x$ を再現している。

さて補間関数 $f(x)$ を用いて以下のようにモデルを構成してみよう。

$$y = Af(x - c) + \epsilon \quad (7.11)$$

$$\epsilon \sim \mathcal{N}(0, \sigma) \quad (7.12)$$

上のモデルに $c = 0.3, \sigma = 0.5, A = 10.0$ を用いて作成した模擬データを用いる。

```

1 xb=np.linspace(0.5,1.7,10)
2 c=0.3
3 sign=0.5
4 A=10.0
5 y=A*np.sin(xb+c)+np.random.normal(0,sign,size=len(xb))

```

したがって HMC-NUTS のためのモデルは以下ようになる。

```

1 import numpyro

```

```

2 import numpyro.distributions as dist
3
4 def model(x,y):
5     c = numpyro.sample('c', dist.Uniform(-0.5, 0.5))
6     A = numpyro.sample('A', dist.Uniform(0.1,100))
7     sigma = numpyro.sample('sigma', dist.Exponential(1.0))
8     mu=A*f(x+c)
9     numpyro.sample('y', dist.Normal(mu, sigma), obs=y)

```

これで、HMC-NUTS をいつものように回すと、ちゃんと推定ができる。今回は図7.8に credible interval だけ示そう。

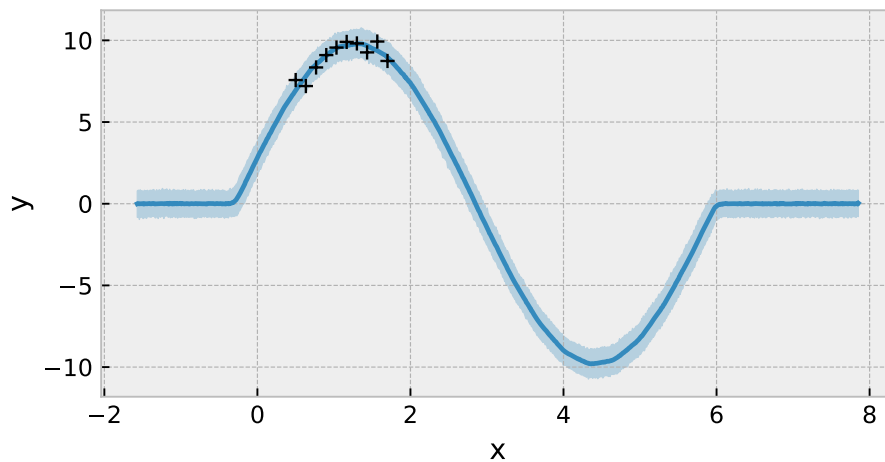


図7.8。

7.1.3 畳み込みの幅を推定する

次の例は、畳み込みの幅を推定する問題である。これは例えばあるシャープなシグナルが、なんらかの過程でなまされてしまったものが観測されるときなどに有用である。典型的なのは観測装置由来のレスポンスである。このような場合 `jax.numpy.convolve` を用いることができる。

レスポンスのカーネルとして次のようなガウス関数型を考えよう。なまされる幅は $\sigma = \text{sigk}$ でパラメトライズされるとする。図7.9に、なまされかたの違いを示す。

さてデータとしてはデルタ関数的なピークが何本かたっているもの (`yp`) を用意して、これをカーネルでなましたものを用いる。HMC 推定としては以下のように構成できる。

```

1 Nk=100
2 xkx=jnp.linspace(-0.1,0.1,Nk)
3

```

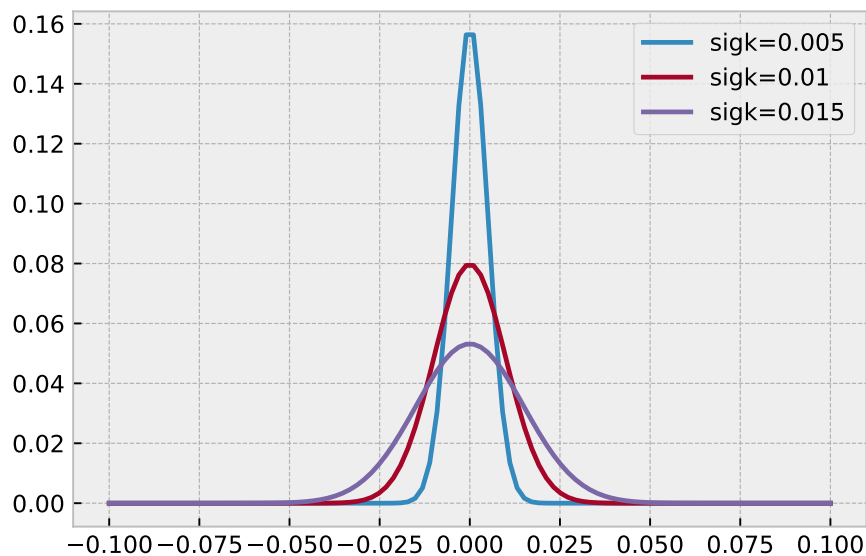


図7.9。

```

4 def model(x,y):
5     sigk = numpyro.sample('sigk', dist.Exponential(1.))
6     sigma = numpyro.sample('sigma', dist.Exponential(1.))
7     mu=jnp.convolve(yp,jnp.exp(-xkx**2/(2.*sigk**2)) \
8                     /jnp.sqrt(2.*jnp.pi)/sigk*(0.2/Nk),"same")
9     numpyro.sample('y', dist.Normal(mu, sigma), obs=y)

```

結果は図7.10の通り。

7.2 カスタマイズした自動微分で numpyro を動かす

Reverse mode

NumPyro の HMC NUTS はデフォルトでは reverse mode を使用する。そこでカスタマイズした自動微分を用いるためには第??章で説明した VJP を定義したものを用いればよい。ここでは $f(x, A) = A \sin x$ の元に、 $y_i = f(x_i, A) + \epsilon$ 、および $\epsilon \sim \mathcal{N}(0, \sigma)$ というモデルで $\mathbf{y} = \{y_i\}$ をフィットしよう。VJP の定義は

```

1 @custom_vjp
2 def h(x,A):
3     return A*jnp.sin(x)
4
5 def h_fwd(x, A):

```

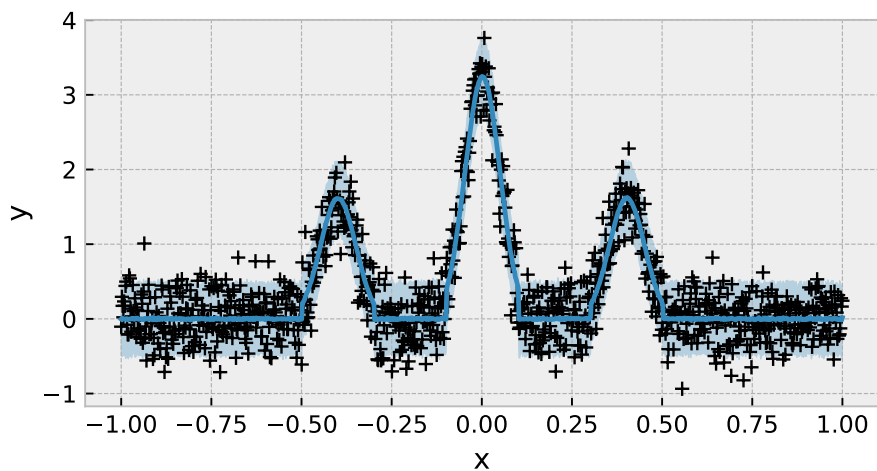


図7.10。

```

6     res = (A*jnp.cos(x), jnp.sin(x))
7     return h(x,A), res
8
9 def h_bwd(res, u):
10     A_cos_x, sin_x = res
11     return (A_cos_x * u, sin_x * u)
12
13 h.defvjp(h_fwd, h_bwd)

```

となる。モデルは

```

1 def model(x,y):
2     sigma = numpyro.sample('sigma', dist.Exponential(1.))
3     x0 = numpyro.sample('x0', dist.Uniform(-1.,1.))
4     A = numpyro.sample('A', dist.Exponential(1.))
5     hv=vmap(h,(0,None),0)
6     mu=hv(x-x0,A)
7     numpyro.sample('y', dist.Normal(mu, sigma), obs=y)

```

とすればよい。ここに `vmap` を用いて x に対してベクトル化していることに注意。HMC-NUTS は

```

1 rng_key = random.PRNGKey(0)
2 rng_key, rng_key_ = random.split(rng_key)
3 num_warmup, num_samples = 1000, 2000
4 kernel = NUTS(model)

```

```
5 | mcmc = MCMC(kernel, num_warmup, num_samples)
6 | mcmc.run(rng_key_, x=x, y=data)
7 | mcmc.print_summary()
```

とすればよい。

Forward mode

NumPyro 0.5.0 より forward mode での HMC NUTS も動作するようになった。

第 8 章

NumPyro で楽しむガウス過程

8.1 ガウス過程

ガウス過程は、多変数正規分布

$$\mathbf{d} \sim \mathcal{N}(\mathbf{0}, \Sigma) \quad (8.1)$$

に従う確率変数 \mathbf{x} の確率過程のことである。いま d_i は時系列だとして、この多変数正規分布の共分散行列の各成分が、

$$\Sigma_{ij} = K_{ij}(a, \tau) = ak(|t_i - t_j|; \tau) \quad (8.2)$$

のように、 t_i と t_j の差の絶対値の関数として与えられれば、相関長 τ のガウス過程が得られる。カーネル関数 $k(t, \tau)$ としてはいろいろなタイプがありうるが、例えば RBF カーネル

$$k_{\text{RBF}}(t; \tau) = \exp\left(-\frac{t^2}{2\tau^2}\right), \quad (8.3)$$

と Matérn 3/2 カーネル

$$k_{\text{M3/2}}(t; \tau) = \left(1 + \frac{\sqrt{3}t}{\tau}\right) e^{-\sqrt{3}t/\tau}. \quad (8.4)$$

などがある。

とりあえず論より実装で、例を作ってみよう。上記カーネルは

```
1 import numpy as np
2 def RBF(t,tau):
3     Dt = t - np.array([t]).T
4     K=np.exp(-(Dt)**2/(tau**2))
5     return K
6
7 def Matern32(t,tau):
8     Dt = t - np.array([t]).T
9     fac=np.sqrt(3.0)*np.abs(Dt)/tau
```



```

10     K=(1.0+fac)*np.exp(-fac)
11     return K

```

のような感じである。さて、`scipy.stats.multivariate_normal` で多変数正規分布がサンプリングできる^{*1}。

```

1  from scipy.stats import multivariate_normal as smn
2  np.random.seed(seed=1) # 乱数初期値
3  N = 101
4  t = np.linspace(0,10,N)
5  ave = np.zeros(N)
6  tau = 0.4
7  a=1.0
8  ave = jnp.zeros(N)
9  cov = a*RBF(t,tau) #cov = Matern32(t,tau) でもよい
10 di = smn(mean=ave ,cov=cov , allow_singular =True).rvs(1).T

```

さて、このように生成されたデータに平均ゼロ、標準偏差 σ のガウスノイズを足しておく。このノイズを便宜的に観測ノイズと呼んでおこう。

```

1  sigma=0.6
2  d=di+np.random.normal(0.0,sigma,len(di))

```

図8.1が生成されたデータである。

さてこのデータを元に τ 、 a 、 σ を推定してみよう。平均ゼロ、標準偏差 σ のガウスノイズの観測ノイズを足したということは確率モデルとしては、

$$\mathbf{d} \sim \mathcal{N}(\mathbf{0}, \Sigma') \quad (8.5)$$

$$\Sigma' = K(a, \tau) + \sigma^2 I \quad (8.6)$$

となる。この Σ' を以下のように実装する。

```

1  def modelcov(t,tau,a,sigma):
2      Dt = t - jnp.array([t]).T
3      K=a*jnp.exp(-(Dt)**2/2/(tau**2))+jnp.eye(N)*sigma**2
4      return K

```

`numpyro` のモデルはこの `modelcov` を用いて

```

1  import jax.numpy as jnp
2  import numpyro
3  import numpyro.distributions as dist

```

^{*1} `numpy.random` や `JAX.random`、および `dist.MultivariateNormal` でもあるのだがとある理由 (9.1章参照) によりここでは `scipy` のものを用いる。

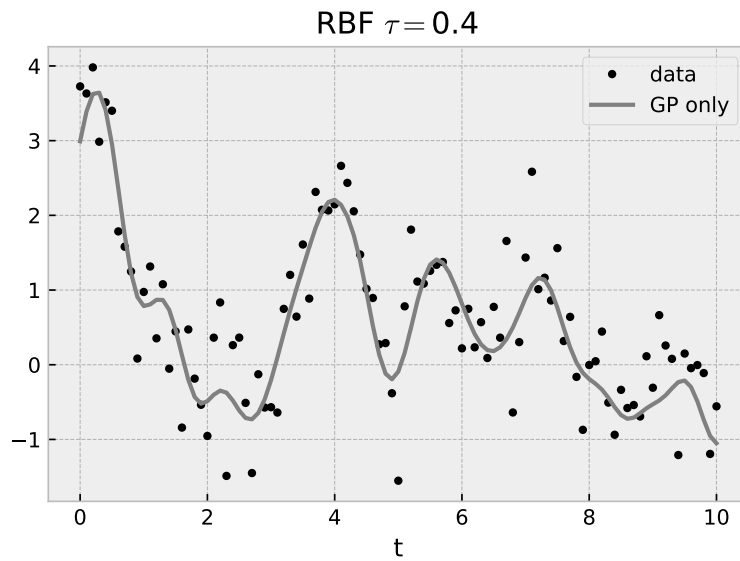


図8.1 実線は σ の観測ノイズなし。

```

4
5 def model(t,y):
6     sigma = numpyro.sample('sigma', dist.Exponential(1.))
7     tau = numpyro.sample('tau', dist.Exponential(1.))
8     a = numpyro.sample('a', dist.Exponential(1.))
9     cov = modelcov(t,tau,a,sigma)
10    numpyro.sample('y', dist.MultivariateNormal(loc=jnp.zeros(N),
11    covariance_matrix=cov), obs=y)

```

となる。ここでは τ と σ の事前分布を Exponential 分布とした。numpyro の HMS-NUTS を以下のように回す。

```

1 from numpyro.infer import MCMC, NUTS
2 from jax import random
3 rng_key = random.PRNGKey(0)
4 rng_key, rng_key_ = random.split(rng_key)
5 num_warmup, num_samples = 1000, 2000
6 # Run NUTS.
7 kernel = NUTS(model)
8 mcmc = MCMC(kernel, num_warmup, num_samples)
9 mcmc.run(rng_key_, t=t, y=d)

```

そして

```

1 import arviz
2 refs={};refs["sigma"]=sigma;refs["tau"]=tau;refs["a"]=a
3 arviz.plot_pair(arviz.from_numpyro(mcmc),kind='kde',
4               divergences=False,marginals=True,reference_values=refs,
5               reference_values_kwargs={'color':"red", "marker":"o", "markersize":12})

```

のように事後分布を可視化したものが図8.2である。ところで、式 (8.5) という観点のモデル化では、 $x = 0$

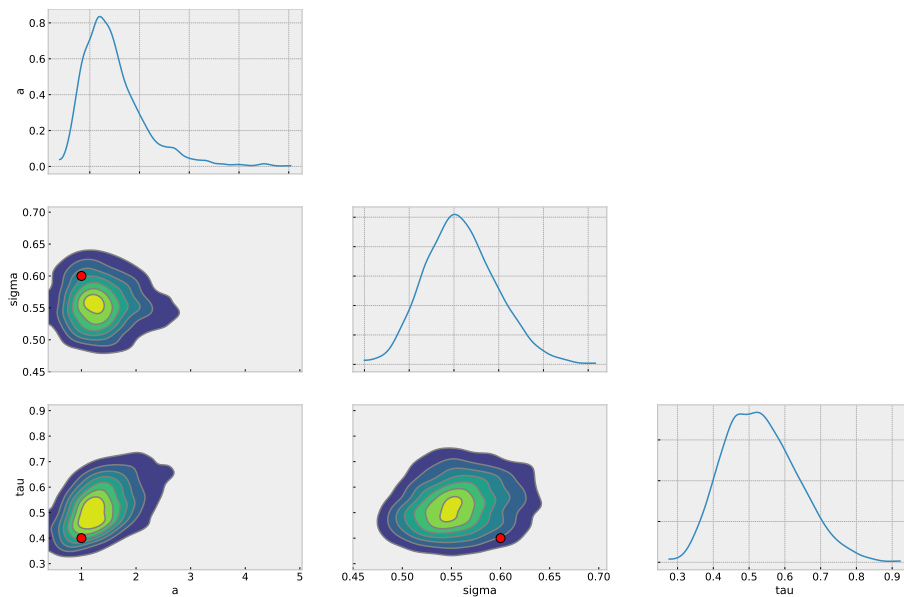


図8.2

がモデル平均である。これは第7章で行ったのと同様に credible interval を計算することでより明確になる。つまり

```

1 from numpyro.infer import Predictive
2 from numpyro.diagnostics import hpdi
3 # 事後分布のサンプリングを取り出す
4 posterior_a = mcmc.get_samples()['a']
5 posterior_tau = mcmc.get_samples()['tau']
6 posterior_sigma = mcmc.get_samples()['sigma']
7 # 予測モデル
8 pred = Predictive(model,{'a':posterior_a,'tau':posterior_tau,
9                          'sigma':posterior_sigma},return_sites=["y"])
10 predictions = pred(rng_key_,t=t,y=None)
11 mean_muy = jnp.mean(predictions["y"], axis=0)
12 hpdi_muy = hpdi(predictions["y"], 0.9)

```

で予測サンプリングから hpdi を計算し、

```
1 fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(6, 3))
2 ax.plot(t,d,".",color="black")
3 ax.plot(t_,mean_muy,color="C0")
4 ax.fill_between(t_, hpdi_muy[0], hpdi_muy[1], alpha=0.3, interpolate=True,color="C0")
5 plt.xlabel("t")
6 plt.ylabel("y")
```

のようにプロットすると、図8.3となり、直感的にわかる。

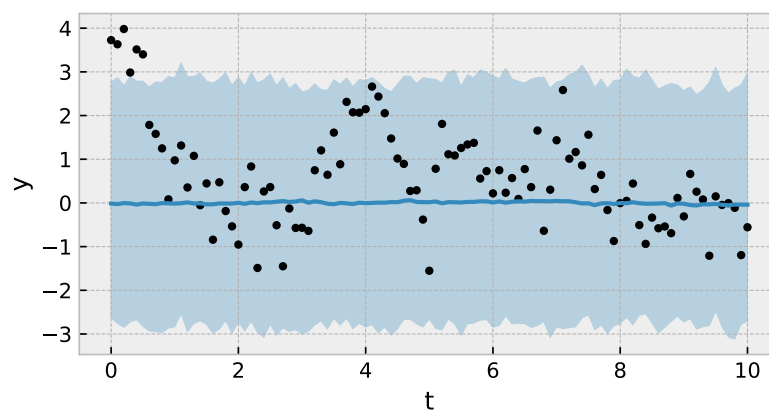


図8.3

モデルパラメタの事前分布としてのガウス過程

上記のガウス過程は、モデルパラメタ \mathbf{m} の事前分布（プライア）として

$$p(\mathbf{m}) = \mathcal{N}(\mathbf{0}, \Sigma) \quad (8.7)$$

と置き、

$$d_i = m_i + \epsilon \quad (8.8)$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2) \quad (8.9)$$

と見ることもできる。ここに ϵ は観測ノイズに対応している。もう少しまどろっこしく書くとモデル \mathbf{g} を恒等変換として、

$$\mathbf{g}(\mathbf{m}) = \mathbf{m} \quad (8.10)$$

$$\mathbf{d} = \mathbf{g}(\mathbf{m}) + \epsilon \quad (8.11)$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2 I) \quad (8.12)$$

のように書ける。この場合、尤度関数が

$$p(\mathbf{d}|\mathbf{m}) = \mathcal{N}(\mathbf{d} - \mathbf{g}(\mathbf{m}), \sigma^2 I) = \mathcal{N}(\mathbf{d} - \mathbf{m}, \sigma^2 I) \quad (8.13)$$

となる。そして、さらにこの事前分布の中にパラメタがあるという構造となっている。つまり

$$\Sigma_{ij} = ak(|t_i - t_j|; \tau) \quad (8.14)$$

としていて、このパラメタ a, τ をハイパーパラメタという。そしてハイパーパラメタの事前分布 (超事前分布またはハイパープライア) をさらに仮定するということに相当する。

以上のモデルを、直接 `numpyro` のモデルとして構成し、 \mathbf{m} の事後分布をサンプリングすることもできるが、ハイパーパラメタを固定した状態では、 \mathbf{m} の事後分布を解析的に書けることが知られている。ここでガウス過程の計算法を紹介しよう。まず多変数正規分布に多変数正規分布をかけても多変数正規分布である。そこで計算しなくてはならないのは \exp のべきの部分のみである。多変数正規分布の「べき」部分は

$$-2 \log \mathcal{N}(\mathbf{m} | \boldsymbol{\mu}, \Sigma) = (\mathbf{m} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{m} - \boldsymbol{\mu}) = \mathbf{m}^\top \Sigma^{-1} \mathbf{m} - 2 \mathbf{m}^\top \Sigma^{-1} \boldsymbol{\mu} + \text{const}. \quad (8.15)$$

となっていることに注意すると、もしある多変数正規分布に従うことがわかっている確率密度分布 $p(\mathbf{m})$ が

$$-2 \log p(\mathbf{m}) = \mathbf{m}^\top P \mathbf{m} - 2 \mathbf{m}^\top \mathbf{q} + \text{const}, \quad (8.16)$$

のように書けたとすると、式 (8.15) と比べることで、この確率密度分布は

$$p(\mathbf{m}) = \mathcal{N}(\mathbf{m} | P^{-1} \mathbf{q}, P^{-1}). \quad (8.17)$$

であることがわかる。

さて、いま事後分布は

$$p(\mathbf{m} | \mathbf{d}) \propto p(\mathbf{d} | \mathbf{m}) p(\mathbf{m}) = \mathcal{N}(\mathbf{d} - \mathbf{m}, \sigma^2 I) \mathcal{N}(\mathbf{0}, \Sigma) \quad (8.18)$$

であるので、この「べき」部分を \mathbf{m} について展開することで、

$$-2 \log [\mathcal{N}(\mathbf{d} - \mathbf{m}, \sigma^2 I) \mathcal{N}(\mathbf{0}, \Sigma)] = \mathbf{m}^\top (\Sigma^{-1} + \sigma^{-2} I) \mathbf{m} - 2 \sigma^{-2} \mathbf{m}^\top \mathbf{d} + \text{const}. \quad (8.19)$$

と書けること、また

$$(\Sigma^{-1} + \sigma^{-2} I)^{-1} = \Sigma (I + \sigma^{-2} \Sigma)^{-1} \quad (8.20)$$

より

$$p(\mathbf{m} | \mathbf{d}) = \mathcal{N}(\Sigma (\sigma^2 I + \Sigma)^{-1} \mathbf{d}, \Sigma (I + \sigma^{-2} \Sigma)^{-1}) \quad (8.21)$$

であることがわかった。

さて、先程、行った HMC-NUTS による τ と σ のサンプリングがこの枠組みで何を意味するか考えよう。これらは今の枠組みではハイパーパラメタとみなせるので、まとめて $\boldsymbol{\theta} = (a, \tau, \sigma)^\top$ をハイパーパラメタとおく。ここからは $\boldsymbol{\theta}$ も確率に入れて考えることにする。さて、今、 \mathbf{m} について周辺化した尤度

$$p(\mathbf{d} | \boldsymbol{\theta}) = \frac{p(\mathbf{d} | \mathbf{m}, \boldsymbol{\theta}) p(\mathbf{m}, \boldsymbol{\theta})}{p(\mathbf{m} | \mathbf{d}, \boldsymbol{\theta})} \quad (8.22)$$

はやはり多変数正規分布となるが、また「べき」部分の \mathbf{d} に関する部分だけ考えることにより

$$-2 \log p(\mathbf{d} | \boldsymbol{\theta}) = -2 \log p(\mathbf{d} | \mathbf{m}, \boldsymbol{\theta}) + 2 \log p(\mathbf{m} | \mathbf{d}, \boldsymbol{\theta}) + \text{const} \quad (8.23)$$

$$= -2 \log \mathcal{N}(\mathbf{d} - \mathbf{m}; \sigma^2 I) + 2 \log \mathcal{N}(\Sigma (\sigma^2 I + \Sigma)^{-1} \mathbf{d}, \Sigma (I + \sigma^{-2} \Sigma)^{-1}) \quad (8.24)$$

$$= \mathbf{d}^\top (\sigma^2 I + \Sigma)^{-1} \mathbf{d} + \text{const}. \quad (8.25)$$

である。よって

$$p(\mathbf{d}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \Sigma + \sigma^2 I) \quad (8.26)$$

となる。これは $\Sigma = K(\tau)$ とすれば、式 (8.5) と同じである。すなわち、(超) 事前分布として $p(\boldsymbol{\theta})$ を与えた時の周辺化事後分布

$$p(\boldsymbol{\theta}|\mathbf{d}) \propto p(\mathbf{d}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (8.27)$$

を HMC-NUTS でサンプリングしていたことになり、図8.3はハイパーパラメタ τ 、 σ の周辺化事後分布を示していることがわかった。つまり $k = 0, \dots, N_s - 1$ に対して

$$\boldsymbol{\theta}_k^\dagger \sim p(\boldsymbol{\theta}|\mathbf{d}) \propto p(\mathbf{d}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (8.28)$$

をサンプリングしたことになる。さて、ではこのサンプルと式 (8.21) を用いて、

$$p(\mathbf{m}, \boldsymbol{\theta}|\mathbf{d}) = p(\mathbf{m}|\boldsymbol{\theta}, \mathbf{d})p(\boldsymbol{\theta}|\mathbf{d}) \quad (8.29)$$

であることから

$$p(\mathbf{m}|\boldsymbol{\theta}_k^\dagger, \mathbf{d}) = \mathcal{N}(\boldsymbol{\mu}_k, K_k) \quad (8.30)$$

$$\boldsymbol{\mu}_k = K(a_k^\dagger, \tau_k^\dagger)((\sigma_k^\dagger)^2 I + K(a_k^\dagger, \tau_k^\dagger))^{-1} \mathbf{d} \quad (8.31)$$

$$K_k = K(a_k^\dagger, \tau_k^\dagger)(I + (\sigma_k^\dagger)^{-2} K(a_k^\dagger, \tau_k^\dagger))^{-1} \quad (8.32)$$

からサンプルした \mathbf{m}_k^\dagger と $\boldsymbol{\theta}_k^\dagger$ のセットは

$$(\mathbf{m}_k^\dagger, \boldsymbol{\theta}_k^\dagger) \sim p(\mathbf{m}, \boldsymbol{\theta}|\mathbf{d}) \quad (8.33)$$

とみなせることがわかる^{*2}。式 (8.31)、(8.32) を以下のように実装する。

```

1 def muGP(tau,a,sigma):
2     cov = a*RBF(t,tau)
3     IKw=sigma**2*np.eye(N)+cov
4     A=scipy.linalg.solve(IKw,d,assume_a="pos")
5     return cov@A
6
7 def covGP(tau,a,sigma):
8     cov = a*RBF(t,tau)
9     IKw=np.eye(N)+cov/sigma**2
10    IKw=scipy.linalg.inv(IKw)
11    return cov@IKw

```

のように定義する。次に式 (8.30) からの \mathbf{m} のサンプリングを

^{*2} ハイパーパラメタの点推定（いわゆる maximum marginal likelihood=maximum evidence）を用いる方法については解説が多くあるが、周辺化事後分布から再サンプリングする議論についてはいまのところ文献が見つからない。ベイズ線形問題+ガウス過程の枠組みで以前、筆者が議論したことがあるので参考までに挙げておく [4]。

```

1 import tqdm
2 Ns=len(posterior_sigma)
3 np.random.seed(seed=1)
4 marr=[]
5 for i in tqdm.tqdm(range(0,Ns)):
6     sigmas=float(posterior_sigma[i])
7     taus=float(posterior_tau[i])
8     a_s=float(posterior_a[i])
9     ave=muGP(taus,a_s,sigmas)
10    cov=covGP(taus,a_s,sigmas)
11    mk = smn(mean=ave ,cov=cov , allow_singular =True).rvs(1).T
12    marr.append(mk)
13 marr=np.array(marr)

```

のように実行し、

```

1 mean_muy = np.mean(marr, axis=0)
2 hpdi_muy = hpdi(marr, 0.9)

```

HPDI を計算することで、図8.4の m の credible interval が求まる。ここで注意点としては、これはモデル m の 90% interval であり、観測ノイズ部分を含んだ d の予測ではないという点である。

```

1 fig=plt.figure(figsize=(6, 3))
2 ax=fig.add_subplot(111)
3 ax.plot(t,d,".",color="black")
4 ax.plot(t, mean_muy,color="C1")
5 ax.fill_between(t, hpdi_muy[0], hpdi_muy[1], alpha=0.3, interpolate=True,color="C1")
6 plt.xlabel("t")
7 plt.ylabel("y")

```

データ点にない位置の予測

さて、一般に $t = t^*$ での予測値はどうなるだろうか？以下では再度ハイパーパラメタを省略して表記する。まず、 m 、 m^* がガウス過程に従うとき

$$p(m^*|m) = \mathcal{N}(K_{\times}^{\top} K^{-1} m, K_{*} - K_{\times}^{\top} K^{-1} K_{\times}) \quad (8.34)$$

である (Appendix A)。ここに

$$K_{ij} = ak(|t_i - t_j|; \tau) \quad (8.35)$$

$$(K_{\times})_{ij} = ak(|t_i - t_j^*|; \tau) \quad (8.36)$$

$$(K_{*})_{ij} = ak(|t_i^* - t_j^*|; \tau) \quad (8.37)$$

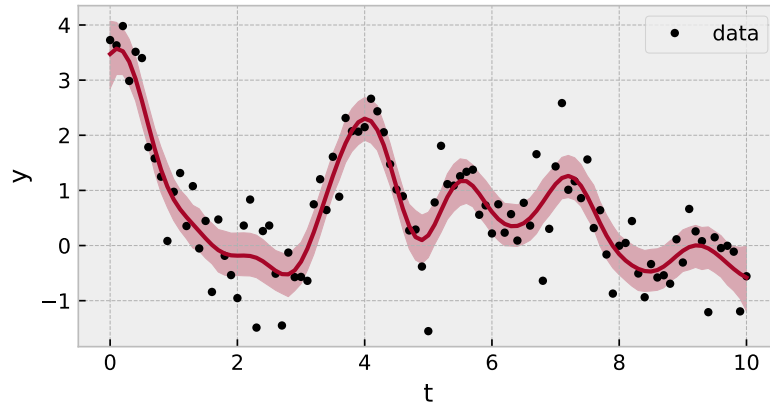


図8.4

同様に

$$p(\mathbf{m}^*|\mathbf{d}) = \mathcal{N}(\mathbf{K}_{\times}^{\top} \mathbf{K}_{\sigma}^{-1} \mathbf{d}, \mathbf{K}_{*} - \mathbf{K}_{\times}^{\top} \mathbf{K}_{\sigma}^{-1} \mathbf{K}_{\times}) \quad (8.38)$$

となる。ここに

$$(\mathbf{K}_{\sigma})_{ij} = ak(|t_i - t_j|; \tau) + \sigma^2 \delta_{i,j} \quad (8.39)$$

($\delta_{i,j}$ はクロネッカーデルタ)。ここで $p(\mathbf{m}^*|\mathbf{d})$ はあくまでモデルパラメタとしてのガウス過程の事後分布なので、観測ノイズの分は考慮されていない。

もし観測ノイズを含んだ予測を行いたいのであれば、

$$\mathbf{d}^* = \mathbf{m}^* + \epsilon \quad (8.40)$$

というモデルに基づき、

$$p(\mathbf{d}^*|\mathbf{d}) = \mathcal{N}(\mathbf{K}_{\times}^{\top} \mathbf{K}_{\sigma}^{-1} \mathbf{d}, \mathbf{K}_{*,\sigma} - \mathbf{K}_{\times}^{\top} \mathbf{K}_{\sigma}^{-1} \mathbf{K}_{\times}) \quad (8.41)$$

となるだろう。ここに

$$(\mathbf{K}_{*,\sigma})_{ij} = ak(|t_i^* - t_j^*|; \tau) + \sigma^2 \delta_{i,j} \quad (8.42)$$

である。

というわけで、ハイパーパラメタ込みのサンプリングで式 (8.38, 8.41) のサンプリングを行い Credible Interval を求める。

```

1 def mucovGPx(t,td,tau,a,sigma):
2     cov = a*RBF(t,tau) + sigma**2*np.eye(N)
3     covx= a*RBFx(t,td,tau)
4     # m の推定の場合
5     covxx = a*RBF(td,tau)
6     # 観測ノイズ込みの推定の場合

```



```

7   covxx = a*RBF(td,tau) + sigma**2*np.eye(N)
8   IKw=cov
9   A=scipy.linalg.solve(IKw,d,assume_a="pos")
10  IKw = scipy.linalg.inv(IKw)
11  return covx@A, covxx - covx@IKw@covx.T # 平均、共分散行列

```

これで同様に HPDI を求めたものが図8.5である。濃い色が m^* の、薄い色が d^* の credible interval である。

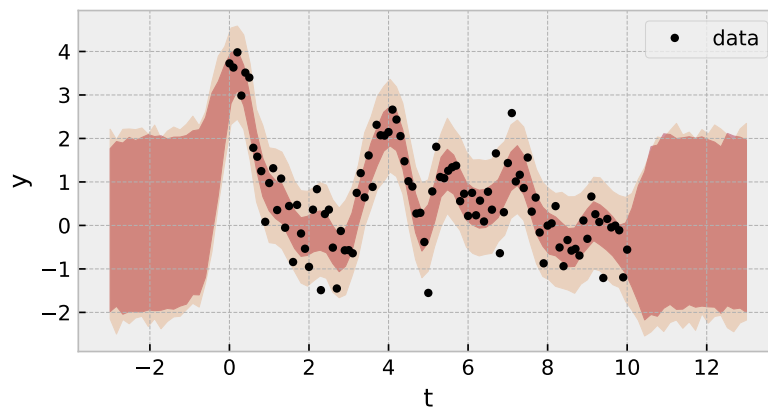


図8.5

8.2 ガウス過程ノイズを含んだモデルフィット

さて、ここまではモデルとしては平均値がゼロのガウス過程により生成されるモデル

$$m \sim \mathcal{N}(0, \Sigma(t)) \quad (8.43)$$

を考えていた。ここでは一般に平均値が t の関数となるガウス過程によりデータが生成される場合を考えよう。すなわち

$$m \sim \mathcal{N}(f(t), \Sigma(t)) \quad (8.44)$$

の場合を考えよう。ここに $f(t)$ は $f(t)$ に対し、 $t = (t_0, t_1, \dots, t_{N-1})$ を要素ごとに適用したベクトル版である。ここでは $f(t)$ として

$$f(t) = ke^{-(t-T_0)^2/2s^2} \sin(2\pi t/P) \quad (8.45)$$

というものを考えてみよう。上では簡単に $f(t)$ と書いたが、パラメタを $\theta = (T_0, k, s, P)$ として $f(t; \theta)$ という表記も併用する。

式 (8.44) の共分散行列を RBF カーネルとして選択すると、例えば、図8.6のようなデータが生成される。ここで $\tau = 3$ であり $f(t)$ より比較的ゆったりとしたトレンドが乗っているのが見て取れるだろう。このようなモデルを HMC フィットするというのは、シグナル $f(t)$ と相関のあるノイズ+観測ノイズをモデル化し、 $f(t)$ の持つパラメタを推定したいときに対応する。

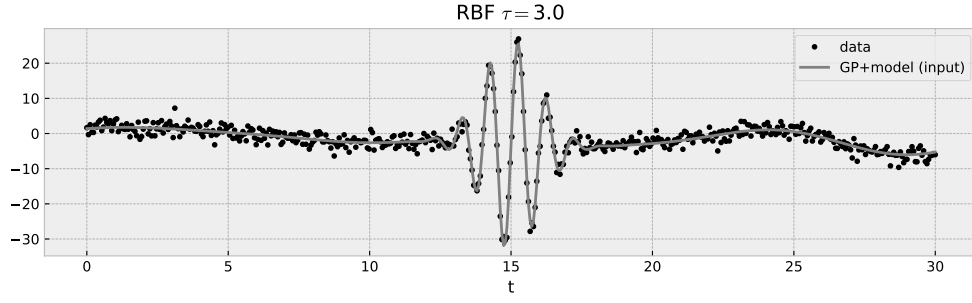


図8.6 実線は σ の観測ノイズなし。

$f(t)$ がわかっている場合、観測ノイズを含まない、および、含んだ予測は、それぞれ、

$$p(\mathbf{m}^*|\mathbf{d}) = \mathcal{N}(\mathbf{f}(t^*) + K_{\times}^{\top} K_{\sigma}^{-1}(\mathbf{d} - \mathbf{f}(t)), K_{*} - K_{\times}^{\top} K_{\sigma}^{-1} K_{\times}) \quad (8.46)$$

$$p(\mathbf{d}^*|\mathbf{d}) = \mathcal{N}(\mathbf{f}(t^*) + K_{\times}^{\top} K_{\sigma}^{-1}(\mathbf{d} - \mathbf{f}(t)), K_{*,\sigma} - K_{\times}^{\top} K_{\sigma}^{-1} K_{\times}) \quad (8.47)$$

となる。 $f(t)$ の持つパラメタ $\theta = (T_0, k, s, P)$ は HMC でサンプリングされるから、例えば後者なら、サンプリングされた各 θ_k^{\dagger} をもちいて

$$\mathbf{d}_k^* \sim \mathcal{N}(\mathbf{f}(t^*; \theta_k^{\dagger}) + K_{\times}^{\top} K_{\sigma}^{-1}(\mathbf{d} - \mathbf{f}(t; \theta_k^{\dagger})), K_{*,\sigma} - K_{\times}^{\top} K_{\sigma}^{-1} K_{\times}) \quad (8.48)$$

をサンプリングすれば、予測のサンプリングができる。というわけで、以下のようにモデル $f(t)$ 込みの GP のフィットを行うことができる。

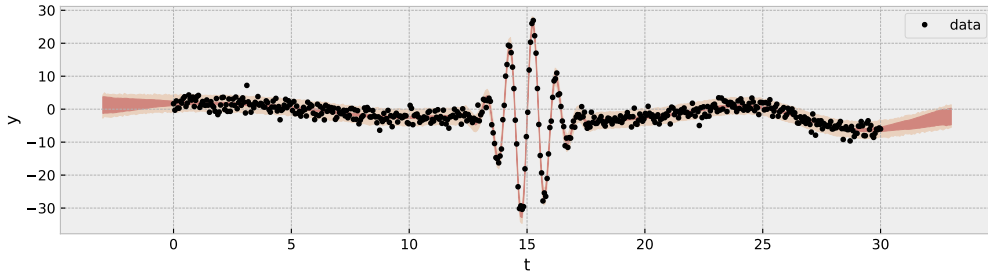


図8.7

第 9 章

NumPyro

NumPyro は PPL なので様々な機能がある。この章では、NumPyro を使っていろいろ遊んでみよう。

9.1 scipy/JAX/NumPyro でサンプリング

scipy/JAX/NumPyro により確率分布からサンプリングを試みよう。例えば、多変数正規分布の場合

```
1 from scipy.stats import multivariate_normal as smn
2 from jax.random import multivariate_normal as jmn
3
4 key = random.PRNGKey(4)
5 N = 1200
6 t = np.linspace(-5,5,N)
7 ave = np.zeros(N)
8 tau = 1.5
9
10 tip=0.1
11 cov = RBF(t,tau)+tip*np.eye(N)
12
13 #scipy
14 ds = smn(mean=ave ,cov=cov , allow_singular =True).rvs(1).T
15 #jax
16 dj=jmn(key,ave,cov)
17 #numpyro
18 mn=dist.MultivariateNormal(loc=ave, covariance_matrix=cov)
19 dn = numpyro.sample('a',mn,rng_key=random.PRNGKey(20))
```

こんな感じで、サンプリングできる。

ここで JAX と NumPyro は default では FP32 を用いていることに注意が必要。これは例えば tip=0.1 を小さくしていくと tip=1.e-6 くらいで JAX と NumPyro では値が nan となってしまう。

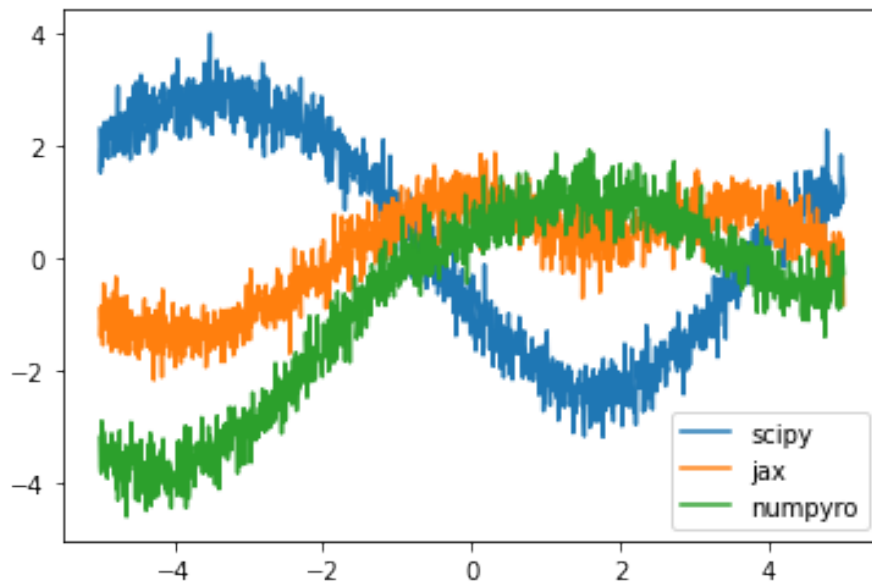


図9.1

```
1 from jax.config import config
2 config.update('jax_enable_x64', True)
```

として FP64 にすれば、JAX/NumPyro でも nan とならない。

key の更新

乱数の種を更新していくには `jax.random.split` を用いる。

```
1 okey=random.PRNGKey(20) #original key
2 for i in range(0,10):
3     okey,key=random.split(okey)
4     dn=numpyro.sample('a',mn,rng_key=key)
```

9.1.1 確率分布の変換

transform

第 10 章

orbax で pytree を保存する

10.1 save/restore

10.1.1 StandardCheckpointer

flax.nnx の state を orbax で保存 (checkpoint) して、restore してみよう。

save は以下のようにする

```
1 import orbax.checkpoint as ocp
2 checkpointer = ocp.StandardCheckpointer()
3 checkpointer.save(ckpt_dir / "state", state)
4 checkpointer.wait_until_finished()
```

最後の行がないと、バッチ実行の場合、保存中に python が終了してコケるので注意が必要である。これは非同期の save を行っているためである。

restore は以下のようにする

```
1 from pathlib import Path
2 ckpt_dir = Path("/home/kawahara/checkpoints")
3 checkpointer = ocp.StandardCheckpointer()
4 restored = checkpointer.restore(ckpt_dir / "state")
```

元の object に restore して再利用したい場合、事前に空箱 (abstract_state) を作成する必要がある。

```
1 from pathlib import Path
2 ckpt_dir = Path("/home/kawahara/checkpoints")
3 checkpointer = ocp.StandardCheckpointer()
4 abstract_model = nnx.eval_shape(
5     lambda: EmuMlpDecoder(rngs=nnx.Rngs(0), grid_length=20000)
6 )
7 graphdef, abstract_state = nnx.split(abstract_model)
8 checkpointer = ocp.StandardCheckpointer()
```

```
9 state_restored = checkpointer.restore(ckpt_dir / "state", abstract_state)
```

10.1.2 Checkpointer

例えば metadata などに加えて複数のデータを保存する場合、Checkpointer を用いて、save 時に args 引数に複数ファイルの保存の仕方をまとめて (Composite) 与える。ここでは state は pytree なので StandardSave を、metadata は JSON で保存したいので JsonSave を使用している。Checkpointer は非同期でないので wait_until_finished() がない。

```
1 metadata = {"version": 1}
2 checkpointer = ocp.Checkpointer(ocp.CompositeCheckpointHandler("state", "metadata"))
3 checkpointer.save(
4     ckpt_dir / "state",
5     args=ocp.args.Composite(
6         state=ocp.args.StandardSave(state), metadata=ocp.args.JsonSave(metadata)
7     ),
8 )
```

restore 時は以下のようにする

```
1 checkpointer = ocp.Checkpointer(ocp.CompositeCheckpointHandler("state", "metadata"))
2 restored = checkpointer.restore(
3     ckpt_dir / "state",
4     args=ocp.args.Composite(
5         state=ocp.args.StandardRestore(abstract_state),
6         metadata=ocp.args.JsonRestore(),
7     ),
8 )
9 state_restored = restored.state
10 metadata = restored.metadata
```

もちろん metadata だけ最初に読むことも可能である。metadata に abstract_state を定義するために必要な情報を格納しておくことで、1 ファイルで restore も可能。

```
1 restored = checkpointer.restore(
2     ckpt_dir / "state",
3     args=ocp.args.Composite(
4         metadata=ocp.args.JsonRestore(),
5     ),
6 )
```

```
7 print(restored.metadata)
```

また `jax.Array` や `numpy.ndarray` を保存するときは `ArraySave` を用いれば良い

```
1 array = jnp.ones(100)
2 checkpointer = ocp.Checkpointer(ocp.CompositeCheckpointHandler("state", "array"))
3 checkpointer.save(
4     ckpt_dir / "state",
5     args=ocp.args.Composite(
6         state=ocp.args.StandardSave(state), metadata=ocp.args.ArraySave(array)
7     ),
8 )
```

restore も同様。

```
1 checkpointer = ocp.Checkpointer(ocp.CompositeCheckpointHandler("state", "array"))
2 restored = checkpointer.restore(
3     ckpt_dir / "state",
4     args=ocp.args.Composite(
5         state=ocp.args.StandardRestore(abstract_state),
6         array=ocp.args.ArrayRestore(),
7     ),
8 )
9 state_restored = restored.state
10 array_restored = restored.array
```

10.1.3 AsyncCheckpointer

`AsyncCheckpointer` は非同期でなので、save 前に終了する条件では `wait_until_finished()` を用いる。

付録 A

Trouble shooting

A.1 JAX 編

JAX の精度が足りてないかもしれない

そんなときは 64bit でチェックだ。

```
1 from jax.config import config
2 config.update("jax_enable_x64", True)
```

微分が nan になってる？

そんなときは以下を書いて、どこで nan が発生しているか調べられる。

```
1 from jax.config import config
2 config.update("jax_debug_nans", True)
```


付録 B

TIPS

B.1 JAX

B.1.1 jax.value_and_grad の has_aux オプション

jax.value_and_grad に has_aux オプションを True とすると、入力関数の帰り値に定義された補助情報を取得できる。

```
1 from jax import value_and_grad
2 import jax.numpy as jnp
3
4 def f(x):
5     value = jnp.sum(x**2)
6     aux = jnp.sum(jnp.abs(x))
7     return value, aux
8
9 x = jnp.array([1.0, 2.0, 3.0])
10 (value, aux), grad = value_and_grad(f, has_aux=True)(x)
11
12 print("value:", value)
13 print("aux:", aux)
14 print("grad:", grad)
```

出力：

```
1 value: 14.0
2 aux: 6.0
3 grad: [2. 4. 6.]
```

補助情報として accuracy や出力などを出力させると便利である。この事情は nnx.value_and_grad でも同じである。

B.2 Arviz

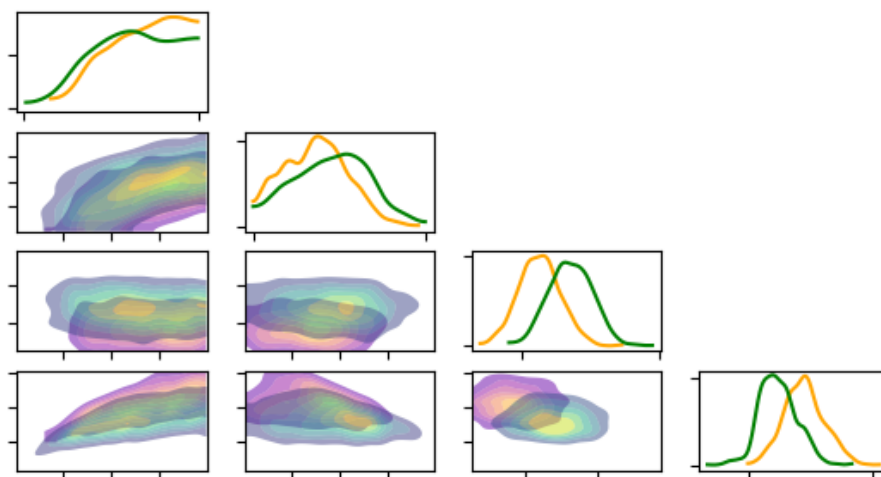
B.2.1 2つのコーナープロットを重ねて描画する

axes を引き継げばよい。

```
1 axes=arviz.plot_pair(p,divergences=False,marginals=True,show=False)
2 axes2=arviz.plot_pair(p2,divergences=False,marginals=True,show=False,ax=axes)
3 plt.show()
```

色等を区別できるようにするためには、各種 kwarg を使い分ければ良い。例えば図B.1の設定は以下のとおりである。

```
1 axes=arviz.plot_pair(p,kind='kde',divergences=False,marginals=True,show=False,\
2     kde_kwargs={"contourf_kwargs":{"alpha":0.5,"cmap":"plasma"},\
3     "contour_kwargs":{"alpha":0}},\
4     marginal_kwargs={"color":"orange"})
5 axes2=arviz.plot_pair(p2,kind='kde',divergences=False,marginals=True,show=False,\
6     ax=axes,\
7     kde_kwargs={"contourf_kwargs":{"alpha":0.5,"cmap":"viridis"},\
8     "contour_kwargs":{"alpha":0}},\
9     marginal_kwargs={"color":"green"})
```



図B.1

付録 C

Appendix

C.1 A. ガウス分布

参考文献

- [1] Mathieu Blondel and Vincent Roulet. The Elements of Differentiable Programming. *arXiv e-prints*, page arXiv:2403.14606, March 2024.
- [2] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- [3] Atilim Gunes Baydin, Barak A. Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *arXiv e-prints*, page arXiv:1502.05767, February 2015.
- [4] Hajime Kawahara and Kento Masuda. Bayesian Dynamic Mapping of an Exo-Earth from Photometric Variability. *ApJ*, 900(1):48, September 2020.
- [5] Radford M. Neal. MCMC using Hamiltonian dynamics. *arXiv e-prints*, page arXiv:1206.1901, June 2012.
- [6] Du Phan, Neeraj Pradhan, and Martin Jankowiak. Composable Effects for Flexible and Accelerated Probabilistic Programming in NumPyro. *arXiv e-prints*, page arXiv:1912.11554, December 2019.
- [7] Robert Edwin Wengert. A simple automatic derivative evaluation program. *Communications of the ACM*, 7(8):463–464, 1964.