

Hessian with Backprop.

Pushpendre Rastogi
January 25, 2021

It is possible to compute the action of the hessian of a function on an input vector. Let $\theta \in \mathbb{R}^d$ and let $f: \mathbb{R}^d \rightarrow \mathbb{R}$.

Let $J = \partial f$ and $H = \partial^2 f$ be the jacobian and hessian function of f . $H(\theta)$ is itself a linear operator and we'll use H_θ to denote this operator.

The trick is that

$$H_\theta \times \nu = \frac{\partial(J_\theta^T \times \nu)}{\partial \theta}.$$

The proof is trivial by just expanding the term on the r.h.s using the product rule and noting that $\frac{\partial \nu}{\partial \theta} = 0$.

This can be implemented in JAX as follows

```
1 from scipy.optimize import (
2     rosen as scipy_rosen, rosen_hess as scipy_rosen_hess)
3 import jax
4 import jax.numpy as jnp
5
6 x = jnp.array([2.2, .21])
7 direction = jnp.array([1.2, 3.1])
8
9 def f(x):
10     return (100.0 * (x[1:] - x[:-1]**2.0)**2.0 + (1 - x[:-1])**2.0)
11     .sum()
12
13 J = jax.grad(f)
14
15 def Jproduct(x, direction):
16     return J(x).dot(direction)
17
18 H_action = jax.grad(Jproduct)
19
20 assert f(x) == scipy_rosen(x), "rosenbrock function impl wrong."
21 gold_action = scipy_rosen_hess(x).dot(direction)
22 trick_action = H_action(x, direction)
23 print(trick_action - gold_action)
24 assert(jnp.linalg.norm(trick_action - gold_action)
25        / jnp.linalg.norm(gold_action) < 1e-6)
26
27 # [-4.8828125e-04 -6.1035156e-05]
```

The second trick is the hutchinson method for extracting the diagonal values of a matrix just by computing the action of the matrix on Rademacher distributed input values. Rademacher R.V. are sampled from a uniform p.m.f. on $\{-1, 1\}^d$. We can see this is in JAX as

```
1 import numpy as np
2 def rademacher(n, d):
3     return np.random.choice([-1, 1], size=n*d).reshape((n,d))
4
5 def approx_diag(x, samples=100):
6     return sum(H_action(x, z) * z for z in rademacher(samples, 2))
7     / samples
8
9 %time print(np.diag(scipy_rosen_hess(x)))
10 %time print(approx_diag(x, 30))
11
12 # [5726.0005  200.      ]
13 # CPU times: user .445 ms, sys: .081 ms, total: .526 ms
14 # Wall time: .464 ms
15
16 # [5667.3335  141.33333]
17 # CPU times: user 1 s, sys: 198 ms, total: 1.2 s
18 # Wall time: 1.03 s
```

See rise.cs.berkeley.edu/projects/adahessian/ and google for adahessian and pyhessian for more details and references.

Now these noisy hessian computations can be used inside optimization methods. The adahessian method is a full-fledged optimization algorithm proposed in "ADAHESIAN: An Adaptive Second Order Optimizer for Machine Learning" based on this idea. See arxiv.org/pdf/2006.00719.pdf

1 Applying Hessian-By-Backprop to Trust Region methods

TODO

Well-known methods such as the trust-region method can also be implemented using the H_action operator.