

Optimization - optax

Lecture 16

Dr. Colin Rundel

SGD Libraries

Most often you will be using the optimizer methods that come with your library of choice, all of the following have their own implementations:

- Tensorflow / Keras
- Torch

Interestingly, JAX does not have builtin support for optimization beyond `jax.scipy.optimize.minimize()` which only supports the BFGS method.

Google previously released `jaxopt` to provide SGD and other optimization methods but this project is now deprecated with the code being merged into DeepMind's Optax.

Optax

Optax is a gradient processing and optimization library for JAX.

Optax is designed to facilitate research by providing building blocks that can be easily recombined in custom ways.

Our goals are to

- Provide simple, well-tested, efficient implementations of core components.
- Improve research productivity by enabling to easily combine low-level ingredients into custom optimizers (or other gradient processing components).
- Accelerate adoption of new ideas by making it easy for anyone to contribute.

We favor focusing on small composable building blocks that can be effectively combined into custom solutions. Others may build upon these basic components in more complicated abstractions. Whenever reasonable, implementations prioritize readability and structuring code to match standard equations, over code reuse.

Same regression example

```
1 from sklearn.datasets import make_regression
2 X, y, coef = make_regression(
3     n_samples=10000, n_features=20, n_informative=4,
4     bias=3, noise=1, random_state=1234, coef=True
5 )
6
7 X = jnp.c_[jnp.ones(len(y)), X]
8 n, k = X.shape
9
10 def lr_loss(beta, X, y):
11     return jnp.sum((y - X @ beta)**2)
```

Optax process

- Construct a `GradientTransformation` object, set optimizer settings

```
1 optimizer = optax.sgd(learning_rate=0.0001); optimizer
```

```
GradientTransformationExtraArgs(init=<function chain.<locals>.init_fn at 0x169936ac0>, update=<function chain.<locals>.update_fn at 0x169936ca0>)
```

- Initialize the optimizer with the initial parameter values

```
1 beta = jnp.zeros(k)
2 opt_state = optimizer.init(beta); opt_state
```

```
(EmptyState(), EmptyState())
```

- Perform iterations

- Calculate the current gradient and update for the optimizer

```
1 f, grad = jax.value_and_grad(lr_loss)(beta, X, y)
2 updates, opt_state = optimizer.update(grad, opt_state); updates, opt_state
```

```
(Array([ 7.1983,  1.8515,  1.1396,  1.7858, -2.8407, -0.1266,
        0.1514, -0.4875, -0.2072, 25.7022, 90.4929,  7.5036,
        0.2313, 123.5414,  1.3136,  2.567 , -0.4262, -1.2996,
        0.5124, -0.2265,  2.3771], dtype=float64), (EmptyState(), EmptyState()))
```

- Apply the update to the parameter

```
1 beta = optax.apply_updates(beta, updates); beta
```

```
Array([ 7.1983,  1.8515,  1.1396,  1.7858, -2.8407, -0.1266,
        0.1514, -0.4875, -0.2072, 25.7022, 90.4929,  7.5036,
        0.2313, 123.5414,  1.3136,  2.567 , -0.4262, -1.2996,
        0.5124, -0.2265,  2.3771], dtype=float64)
```

Basic Example - GD

Implementation

Results

```
1 optimizer = optax.sgd(learning_rate=0.00001)
2
3 beta = jnp.zeros(k)
4 opt_state = optimizer.init(beta)
5
6 gd_loss = []
7 for iter in range(50):
8     f, grad = jax.value_and_grad(lr_loss)(beta, X, y)
9     updates, opt_state = optimizer.update(grad, opt_state)
10    beta = optax.apply_updates(beta, updates)
11    gd_loss.append(f)
12
13 beta
```

```
Array([ 3.0082,  0.009 ,  0.0003,  0.0023,  0.0035,  0.0033,  0.0261,
        -0.0006,  0.0005, 12.2771, 44.4933,  3.6423,  0.0168, 61.3929,
        -0.0011, -0.0054,  0.0139, -0.0094, -0.0054,  0.0023,  0.0219],      dtype=float64)
```

Basic Optax Example - Adam

Implementation

Results

```
1 optimizer = optax.adam(learning_rate=1, b1=0.9, b2=0.999, eps=1e-8)
2
3 beta = jnp.zeros(k)
4 opt_state = optimizer.init(beta)
5
6 adam_loss = []
7 for iter in range(50):
8     f, grad = jax.value_and_grad(lr_loss)(beta, X, y)
9     updates, opt_state = optimizer.update(grad, opt_state)
10    beta = optax.apply_updates(beta, updates)
11    adam_loss.append(f)
12
13 beta
```

```
Array([ 3.3313,  0.229 ,  0.0878,  0.1995, -0.2172, -0.0805,  0.0505,
        -0.033 , -0.2049, 11.6465, 39.4043,  3.8843,  0.1239, 43.1531,
         0.2196,  0.3237, -0.0915, -0.2611,  0.0874,  0.1323,  0.3156],      dtype=float64)
```

A bit more on learning rate and batch size

Optax and mini batches

```
1 def optax_optimize(params, X, y, loss_fn, optimizer, steps=50, batch_size=1, seed=1234):
2     n, k = X.shape
3     res = {"loss": [], "epoch": np.linspace(0, steps, int(steps*(n/batch_size) + 1))}
4
5     opt_state = optimizer.init(params)
6     grad_fn = jax.grad(loss_fn)
7
8     rng = np.random.default_rng(seed)
9     batches = np.array(range(n))
10    rng.shuffle(batches)
11
12    for iter in range(steps):
13        for batch in batches.reshape(-1, batch_size):
14            res["loss"].append(loss_fn(params, X, y).item())
15            grad = grad_fn(params, X[batch,:], y[batch])
16            updates, opt_state = optimizer.update(grad, opt_state)
17            params = optax.apply_updates(params, updates)
18
19    res["params"] = params
20    res["loss"].append(loss_fn(params, X, y).item())
21
22    return(res)
```

Fitting - SGD - Fixed LR (small)

Implementation

Results

```
1 batch_sizes = [10, 100, 1000, 10000]
2 lrs = [0.00001] * 4
3
4 sgd = {
5     batch_size: optax_optimize(
6         params=jnp.zeros(k), X=X, y=y, loss_fn=lr_loss,
7         optimizer=optax.sgd(learning_rate=lr),
8         steps=30, batch_size=batch_size, seed=1234
9     )
10     for batch_size, lr in zip(batch_sizes, lrs)
11 }
```

Fitting - SGD - Adjusted LR

Implementation

Full

Zoom

```
1 batch_sizes = [10, 100, 1000, 10000]
2 lrs = [0.005, 0.001, 0.0001, 0.00001]
3
4 sgd = {
5     batch_size: optax_optimize(
6         params=jnp.zeros(k), X=X, y=y, loss_fn=lr_loss,
7         optimizer=optax.sgd(learning_rate=lr),
8         steps=30, batch_size=batch_size, seed=1234
9     )
10     for batch_size, lr in zip(batch_sizes, lrs)
11 }
```

Fitting - SGD - Fixed LR, Small batch size

Implementation

Full

Zoom

```
1 batch_sizes = [10, 25, 50, 100]
2 lrs = [0.001] * 4
3
4 sgd = {
5     batch_size: optax_optimize(
6         params=jnp.zeros(k), X=X, y=y, loss_fn=lr_loss,
7         optimizer=optax.sgd(learning_rate=lr),
8         steps=2, batch_size=batch_size, seed=1234
9     )
10     for batch_size, lr in zip(batch_sizes, lrs)
11 }
```

Runtime per epoch

Implementation	Runtimes	Scaled
<pre>1 batch_sizes = [10, 50, 100, 10000] 2 lrs = [0.001] * 4 3 4 sgd_runtime = { 5 batch_size: timeit.Timer(lambda: 6 optax_optimize(7 params=jnp.zeros(k), X=X, y=y, loss_fn=lr_loss, 8 optimizer=optax.sgd(learning_rate=lr), 9 steps=1, batch_size=batch_size, seed=1234 10) 11).repeat(5,1) 12 for batch_size, lr in zip(batch_sizes, lrs) 13 }</pre>		

Some lessons / comments

- Batch size determines both training time and computing resources
- Generally there will be an inverse relationship between learning rate and batch size
- Most optimizer hyperparameters are sensitive to batch size
- For really large models batches are a necessity and sizing is often determined by resource / memory constraints

Adam

Adam - Fixed LR

Implementation

Results

```
1 batch_sizes = [10, 25, 50, 100]
2 lrs = [1]*4
3
4 adam = {
5     batch_size: optax_optimize(
6         params=jnp.zeros(k), X=X, y=y, loss_fn=lr_loss,
7         optimizer=optax.adam(learning_rate=lr, b1=0.9, b2=0.999, eps=1e-8),
8         steps=2, batch_size=batch_size, seed=1234
9     )
10 }
11 }
```


Adam - Smaller Fixed LR

Implementation

Results

```
1 batch_sizes = [10, 25, 50, 100]
2 lrs = [0.1]*4
3
4 adam = {
5     batch_size: optax_optimize(
6         params=jnp.zeros(k), X=X, y=y, loss_fn=lr_loss,
7         optimizer=optax.adam(learning_rate=lr, b1=0.9, b2=0.999, eps=1e-8),
8         steps=10, batch_size=batch_size, seed=1234
9     )
10     for batch_size, lr in zip(batch_sizes, lrs)
11 }
```

Learning rate schedules

As mentioned last time, most gradient based methods are not guaranteed to converge unless their learning rates decay as a function of step number.

Optax supports a **large number** of pre-built learning rate schedules which can be passed into any of its optimizers instead of a fixed floating point value.

```
1 schedule = optax.linear_schedule(  
2     init_value=1., end_value=0., transition_steps=5  
3 )  
4  
5 [schedule(step).item() for step in range(6)]
```

```
[1.0, 0.8, 0.6, 0.4, 0.19999999999999996, 0.0]
```

Adam w/ Exp Decay

Implementation

Results

```
1 batch_sizes = [10, 25, 50, 100]
2
3 adam = {
4     batch_size: optax.optimize(
5         params=jnp.zeros(k), X=X, y=y, loss_fn=lr_loss,
6         optimizer=optax.adam(
7             learning_rate=optax.schedules.exponential_decay(
8                 init_value=1,
9                 transition_steps=100,
10                decay_rate=0.9
11            ),
12            b1=0.9, b2=0.999, eps=1e-8
13        ),
14        steps=2, batch_size=batch_size, seed=1234
15    )
16    for batch_size in batch_sizes
17 }
```

Runtime per epoch

Implementation	Runtimes	Scaled
<pre>1 batch_sizes = [10, 25, 50, 100] 2 3 adam_runtime = { 4 batch_size: timeit.Timer(lambda: 5 optax_optimize(6 params=jnp.zeros(k), X=X, y=y, loss_fn=lr_loss, 7 optimizer=optax.adam(8 learning_rate=optax.schedules.exponential_decay(9 init_value=1, 10 transition_steps=100, 11 decay_rate=0.9 12), 13 b1=0.9, b2=0.999, eps=1e-8 14), 15 steps=1, batch_size=batch_size, seed=1234 16) 17).repeat(5,1) 18 for batch_size in batch_sizes 19 }</pre>		

Some advice ...

The following is from Google Research's [Tuning Playbook](#):

- No optimizer is the “best” across all types of machine learning problems and model architectures. Even just comparing the performance of optimizers is a difficult task. 🤖
- We recommend sticking with well-established, popular optimizers, especially when starting a new project.
 - Ideally, choose the most popular optimizer used for the same type of problem.
- Be prepared to give attention to *all* hyperparameters of the chosen optimizer.
 - Optimizers with more hyperparameters may require more tuning effort to find the best configuration.
 - This is particularly relevant in the beginning stages of a project when we are trying to find the best values of various other hyperparameters (e.g. architecture hyperparameters) while treating optimizer hyperparameters as nuisance parameters.
 - It may be preferable to start with a simpler optimizer (e.g. SGD with fixed momentum or Adam with fixed ϵ , β_1 , and β_2) in the initial stages of the project and switch to a more general optimizer later.
- Well-established optimizers that we like include (but are not limited to):
 - SGD with momentum (we like the Nesterov variant)
 - Adam and NAdam, which are more general than SGD with momentum. Note that Adam has 4 tunable hyperparameters and they can all matter!

Optimization in R

Basic optimization

The equivalent of `scipy`'s `optimize.minimize()` for unconstrained continuous optimization problems in R is `stats::optim()` - there is nearly a 1-to-1 correspondence between the two functions and the available optimizers.

```
1 optim(par, fn, gr = NULL, ...,
2       method = c("Nelder-Mead", "BFGS", "CG", "L-BFGS-B", "SANN",
3                  "Brent"),
4       lower = -Inf, upper = Inf,
5       control = list(), hessian = FALSE)
```

The only missing method from `scipy` is `Newton-CG` and there is the addition of the `SANN` method which is a variant of simulated annealing and does not require gradient information. However, it is slow and very sensitive to its control parameters and is not considered a general-purpose method.

All other tuning knobs are hidden in `control` - see the documentation for details. Most important options include: `maxit`, `abstol`, and `reltol`.

Return values

`optim()` returns a list of results, most of which are expected: `par` the minimizer, `value` objective function at `par`, `counts` the number of function and gradient evaluations.

“Success” of the optimization is reported by `convergence` which is a little bit weird (think unix exit codes):

- `0` - indicates successful convergence based on the criteria specified by `control`
- `1` - indicates failure due to reaching the `maxit` limit
- Any other number indicates a special case depending on the method, check `message`

Usage

```
1 ## Rosenbrock Banana function
2 f = function(x) {
3   100 * (x[2] - x[1] * x[1]) ^ 2 + (1 - x[1]) ^ 2
4 }
5 grad = function(x) {
6   c(-400 * x[1] * (x[2] - x[1] * x[1]) - 2 * (1 - x[1]),
7     200 * (x[2] - x[1] * x[1]))
8 }
9 x0 = c(-1.2, 1)
```

```
1 optim(x0, f, grad, method = "BFGS")
```

```
$par
[1] 1 1

$value
[1] 9.594956e-18

$counts
function gradient
      110      43

$convergence
[1] 0

$message
NULL
```

```
1 optim(x0, f, grad, method = "CG")
```

```
$par
[1] -0.7648373  0.5927588

$value
[1] 3.106579

$counts
function gradient
      402      101

$convergence
[1] 1

$message
NULL
```

SGD related methods

For any of these algorithms you will generally be depending on the underlying modeling library to make them available to you, for example:

- Keras [optimizers](#) implemented
- Torch [optimizers](#)

Details are library dependent.

optimx

optimx is an R package that extends and enhances the `optim()` function of base R, in particular by unifying the call to many solvers.

Makes a variety of solvers from different packages available with a unified calling framework.

Packages include: `pracma`, `minqa`, `dfoptim`, `lbfgs`, `lbfgsb3c`, `marqLevAlg`, `nloptr`, `dfoptim`, `BB`, `subplex`, and `ucminf`

nloptr

Wrapper around the [NLopt](#) library (which also has a Python interface).

- Provides a large number of global and local solvers (including everything available in optim)
- Provides more robust support for constrained optimization problems

Usage

```
1 ## Rosenbrock Banana function
2 f = function(x) {
3     100 * (x[2] - x[1] * x[1]) ^ 2 + (1 - x[1])
4 }
5
6 grad = function(x) {
7     c(-400 * x[1] * (x[2] - x[1] * x[1]) - 2 * x[1],
8       200 * (x[2] - x[1] * x[1]))
9 }
10
11 x0 = c(-1.2, 1)
```

```
1 nloptr::nloptr(
2     x0 = x0,
3     eval_f = f, eval_grad_f = grad,
4     opts = list(
5         "algorithm" = "NLOPT_LD_LBFGS",
6         "xtol_rel" = 1.0e-8
7     )
8 )
```

Call:

```
nloptr::nloptr(x0 = x0, eval_f = f, eval_grad_f = grad, opts = list(algorithm = "NLOPT_LD_LBFGS", xtol_rel = 1e-08))
```

Minimization using NLopt version 2.7.1

NLopt solver status: 1 (NLOPT_SUCCESS: Generic success return value.)

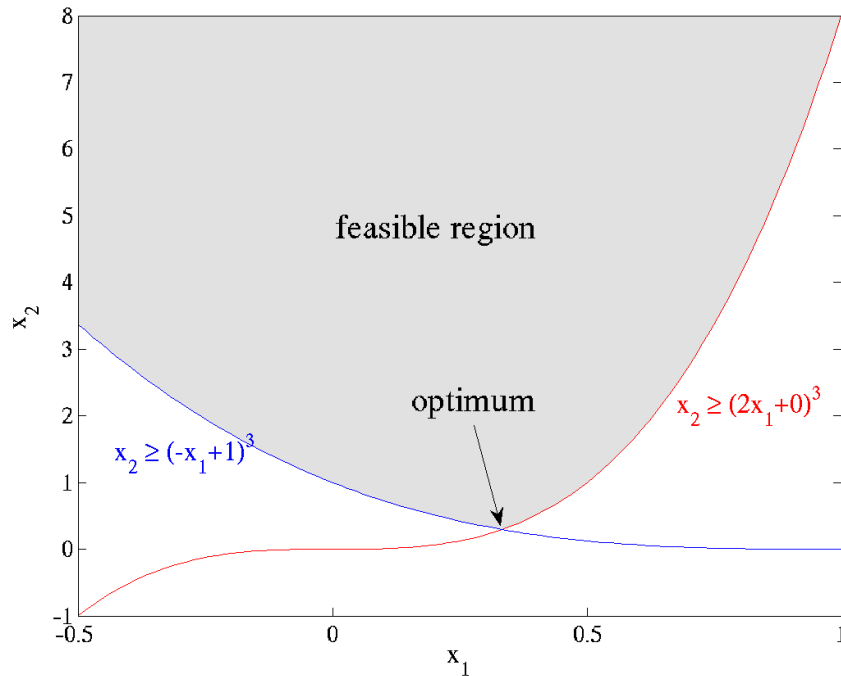
Number of Iterations.....: 56

Termination conditions: xtol_rel: 1e-08

Number of inequality constraints: 0

Number of equality constraints: 0

Constrained Example



$$\begin{aligned} \min_{x \in \mathbb{R}^n} & \sqrt{x_2} \\ \text{s.t.} & \quad x_2 \geq 0 \\ & \quad (a_1 x_1 + b_1)^3 - x_2 \leq 0 \\ & \quad (a_2 x_1 + b_2)^3 - x_2 \leq 0 \end{aligned}$$

where $a_1 = 2$, $b_1 = 0$, $a_2 = -1$, and $b_2 = 1$.

Implementation

```
1 # Objective function & gradient
2 f = function(x, a, b) {
3   sqrt(x[2])
4 }
5 grad_f = function(x, a, b) {
6   c(0, 0.5 / sqrt(x[2]))
7 }
8
9 # Constraint function
10 g = function(x, a, b) {
11   (a * x[1] + b) ^ 3 - x[2]
12 }
13
14 # Jacobian of constraint
15 jac_g = function(x, a, b) {
16   rbind(
17     c(3 * a[1] * (a[1] * x[1] + b[1]) ^ 2,
18       c(3 * a[2] * (a[2] * x[1] + b[2]) ^ 2,
19     )
20   }
21
22 a = c(2, -1)
23 b = c(0, 1)
24
```

Call:

```
nloptr::nloptr(x0 = c(1.234, 5.678), eval_f =
f, eval_grad_f = grad_f,
              lb = c(-Inf, 0), ub = c(Inf, Inf),
eval_g_ineq = g, eval_jac_g_ineq = jac_g,
              opts = list(algorithm = "NLOPT_LD_MMA",
xtol_rel = 1e-08),      a = a, b = b)
```

Minimization using NLopt version 2.7.1

NLopt solver status: 4 (NLOPT_XTOL_REACHED:
Optimization stopped because
xtol_rel or xtol_abs (above) was reached.)

Number of Iterations : 19