# Rapid Optimization Library:
# Machine Learning with PyROL

Brian Chen, Denis Ridzal

Trilinos User-Developer Group Meeting

October 22, 2024

- ROL (as in *rock and **roll***) is a high-performance Trilinos library for **numerical optimization**.

- Brings an extensive collection of modern optimization algorithms to **any application**.

- The programming interface supports **any computational hardware**, including heterogeneous many-core systems with digital and analog accelerators.

- Used successfully in optimal control, optimal design, inverse problems, image processing and mesh optimization, at **extreme problem scales** (very small and very large).

- Application areas including geophysics, structural dynamics, fluid dynamics, electromagnetics, quantum computing, hypersonics and geospatial imaging.

**R L**

RAPID OPTIMIZATION LIBRARY

*Numerical optimization made practical:*
*Any application, any hardware, any problem size.*

- Modern optimization algorithms: **inexact, adaptive, stochastic/risk-aware, nonsmooth**.

- Special programming interfaces for simulation-based optimization: **SimOpt**.

- Toolboxes: **OED** for optimal experimental design and **PDE-OPT** for PDE-constrained optimization.

`rol.sandia.gov`

# Mathematical Formalism

- ROL solves **smooth nonlinear nonconvex optimization** problems

$$\underset{x}{\text{minimize}} \ J(x) \ \text{subject to} \ \begin{cases} c(x) = 0 \\ \ell \leq x \leq u \\ Ax = b \end{cases}$$

where $J : \mathcal{X} \to \mathbb{R}$, $c : \mathcal{X} \to \mathcal{C}$ and $A : \mathcal{X} \to \mathcal{D}$, and $\mathcal{X}$, $\mathcal{C}$ and $\mathcal{D}$ are vector spaces.

- ROL additionally solves **stochastic optimization** problems with random inputs $\xi$,

$$\underset{x}{\text{minimize}} \ \mathcal{R}\left[J(x; \xi)\right] \ \text{etc.} \ \text{with} \ J = J(x; \xi) \ \text{and} \ c = c(x; \xi)$$

where $\mathcal{R}$ is a risk measure, for instance the expectation, $\mathcal{R}\left[J(x; \xi)\right] = \mathbb{E}\left[J(x; \xi)\right]$.

- Finally, ROL solves **nonsmooth optimization** problems

$$\underset{x}{\text{minimize}} \ J(x) + \phi(x)$$

where $\phi : \mathcal{X} \to \mathbb{R}$ is nonsmooth and convex, for instance an $\ell_1$ regularizer.

# A Growing Collection of Algorithms

## Type U "Unconstrained"

minimize $J(x)$
$x$

subject to
$$\begin{cases} \\ Ax = b \end{cases}$$

*Methods*:

- trust region and line search globalization
- gradient descent
- quasi-Newton and inexact Newton
- nonlinear conjugate gradient (CG)
- Cauchy point, dogleg
- Steihaug-Toint truncated CG

## Type B "Bound Constrained"

minimize $J(x)$
$x$

subject to
$$\begin{cases} \ell \le x \le u \\ Ax = b \end{cases}$$

*Methods*:

- projected gradient with line search
- projected Newton with line search
- primal-dual active set
- Lin-Moré trust region
- Kelley-Sachs trust region
- spectral projected gradient (SPG)

## Type E "Equality Constrained"

minimize $J(x)$
$x$

subject to
$$\begin{cases} c(x) = 0 \\ Ax = b \end{cases}$$

*Methods*:

- composite-step sequential quadratic programming (SQP)
- augmented Lagrangian (AL)

## Type G "General Constraints"

minimize $J(x)$
$x$

subject to
$$\begin{cases} c(x) = 0 \\ \ell \le x \le u \\ Ax = b \end{cases}$$

*Methods*:

- AL for equalities and TypeB for bounds
- AL for bounds and TypeE for equalities
- primal interior point
- Moreau-Yosida
- stabilized linearly constrained Lagrangian (LCL)

## Type P "Proximable"

minimize $J(x) + \phi(x)$
$x$

where
$J$ smooth+nonconvex
$\phi$ nonsmooth+convex

*Methods*:

- nonsmooth inexact trust-region methods
- proximal gradient
- spectral proximal gradient
- inexact proximal Newton
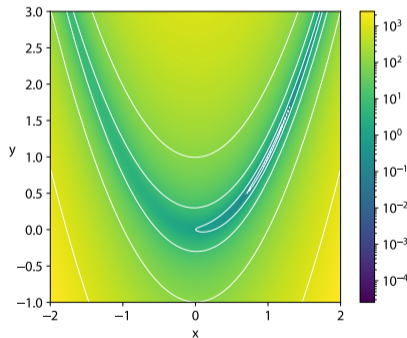
# Motivating PyROL: Rosenbrock Function

$$f(x_1, \ldots, x_n) = \sum_{i=1}^{N/2} \left[ 100(x_{2i-1}^2 - x_{2i})^2 + (x_{2i-1} - 1)^2 \right]$$

- The Rosenbrock function is a nonconvex function that can be difficult to optimize.

- We highlight ROL's capability to optimize over millions of parameters, $N = 100$ million.

- We will use ROL's **Python** interface. Why?

  Python has had wide reach:

  - efficient and easy-to-learn scripting and programming tool;
  - convenient for teaching and prototyping;
  - "glue code" for large projects;
  - popular, including in AI, e.g., JAX, PyTorch.

Rosenbrock function in $2D$ [1]

---
[1] https://commons.wikimedia.org/w/index.php?curid=114931732

$$\underset{x}{\text{minimize}} \quad \sum_{i=1}^{N/2} \left[ 100(x_{2i-1}^2 - x_{2i})^2 + (x_{2i-1} - 1)^2 \right]$$

```
1   class RosenbrockObjective(Objective):
2       def __init__(self):
3           super().__init__()
4           self.alpha = 100
5       def value(self, x, tol):
6           return self.alpha*(x[0]**2-x[1])**2 + (x[0]-1)**2
7       def gradient(self, g, x, tol):
8           g[0] =  4*self.alpha*(x[0]**2-x[1])*x[0] + 2*(x[0]-1)
9           g[1] = -2*self.alpha*(x[0]**2-x[1])
10      def hessVec(self, hv, v, x, tol):
11          h11 = 12*self.alpha*x[0]**2 - 4*self.alpha*x[1] + 2
12          h12 = -4*self.alpha*x[0]
13          h22 =  2*self.alpha
14          hv[0] = h11*v[0] + h12*v[1]
15          hv[1] = h12*v[0] + h22*v[1]
```

```
1   def main():
2       # Set up
3       x = NumPyVector(np.array([-3.,-4.]))
4       objective = RosenbrockObjective()
5       problem = Problem(objective, x)
6       problem.check(True)  # Optional
7       parameters = ParameterList()
8       # Solve
9       solver = Solver(problem, parameters)
10      solver.solve(getCout())
```

$N = 2$

- ROL is hardware agnostic.
- You can run ROL on personal computers (in serial and MPI parallel), on GPUs, and on supercomputers by inheriting from `ROL::Vector`.

```python
1  class NumPyVector(PythonVector):
2      # ...
3      def axpy(self, alpha, other):
4          self.array += alpha*other.array
5      def dot(self, other):
6          return np.vdot(self.array,
7                         other.array)
8      def scale(self, alpha):
9          self.array *= alpha
```

```python
1   class TensorVector(PythonVector):
2       # ...
3       @torch.no_grad()
4       def axpy(self, alpha, other):
5           self.tensor.add_(other.tensor, alpha=alpha)
6       @torch.no_grad()
7       def dot(self, other):
8           ans = torch.sum(torch.mul(self.tensor,
9                                     other.tensor))
10          return ans.item()
11      @torch.no_grad()
12      def scale(self, alpha):
13          self.tensor.mul_(alpha)
```

# PyTorch Automatic Differentiation

```
1   class TorchObjective(Objective):
2       def __init__(self):
3           super().__init__()
4           self.torch_gradient = torch.func.grad(self.torch_value)
5       def torch_value(self, x):
6           # Returns a scalar torch Tensor
7           raise NotImplementedError
8       def value(self, x, tol):
9           return self.torch_value(x.torch_object).item()
10      def gradient(self, g, x, tol):
11          ans = self.torch_gradient(x.torch_object)
12          g.torch_object = ans
13      def hessVec(self, hv, v, x, tol):
14          input = torch.func.grad(self.torch_value)
15          _, ans = self._forward_over_reverse(input, x.torch_object, v.torch_object)
16          hv.torch_object = ans
17      def _forward_over_reverse(self, input, x, v):
18          # https://github.com/google/jax/blob/main/docs/notebooks/autodiff_cookbook.ipynb
19          return torch.func.jvp(input, (x,), (v,))
```

# PyTorch Rosenbrock Example

$$\underset{x}{\text{minimize}} \quad \sum_{i=1}^{N/2} \left[ 100(x_{2i-1}^2 - x_{2i})^2 + (x_{2i-1} - 1)^2 \right]$$

```
1   class RosenbrockObjective(TorchObjective):
2       def torch_value(self, x):
3           return torch.sum(100*(x[::2]**2 - x[1::2])**2 + (x[::2] - 1)**2)
```

PyTorch automatically differentiates this problem, and we can also leverage its GPU dispatch:

```
1   device = torch.device('cuda')
2   x = torch.empty(N, requires_grad=False, device=device)
3   x = TensorVector(x)
```

When $N = 10^8$, the problem takes:
$\sim$ 466 seconds to solve with an Intel Core i9 CPU; and
$\sim$ 7 seconds with an NVIDIA GPU.

Denis Ridzal                                                                                    ROL: Machine Learning with PyROL

# Neural Network Example

Set up a neural network in PyTorch.

```
1   class ConvolutionalNet(nn.Module):
2       def __init__(self):
3           super(ConvolutionalNet, self).__init__()
4           self.conv1 = nn.Conv2d(1, 8, 3, 1)
5           self.conv2 = nn.Conv2d(8, 8, 3, 1)
6           self.fc1 = nn.Linear(1152, 128)
7           self.fc2 = nn.Linear(128, 10)
8
9       def forward(self, x):
10          x = .....
11          output = F.log_softmax(x, dim=1)
12          return output
```

Wrap the parameters using `TensorDictVector` and pass the model into a `TorchObjective`.

```
1   def main(data, model, loss_fcn):
2       x = TensorDictVector(model.state_dict())
3       objective = LeastSquaresObjective(data, model)
4       g = TensorDictVector(
5           copy.deepcopy(model.state_dict())
6       )
7
8       stream = getCout()
9       problem = Problem(objective, x, g)
10      problem.checkDerivatives(True, stream)
11
12      params = build_parameter_list(iteration_limit)
13      solver = Solver(problem, params)
14      solver.solve(stream)
15
16      return solver, x
```
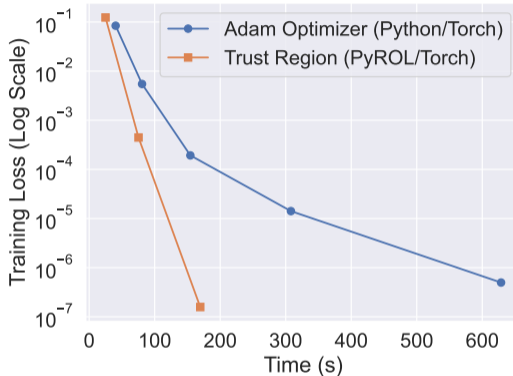
# Example: MNIST Classification

- <u>Goal</u>: Minimize classification error on MNIST using a convolutional neural network (CNN).

- <u>Optimization Parameters</u>: Around $100,000$ parameters of the CNN.



Images from the MNIST handwritten digits dataset.



MNIST: Comparison of training time and final training errors, between Adam (implemented using PyTorch) versus ROL's trust region (implemented using PyROL and PyTorch).