

Learning with differentiable and amortized optimization

Brandon Amos • Meta AI (FAIR) NYC



<http://github.com/bamos/presentations>

Optimization is crucial technology

Optimization is a **modeling** and **decision-making** paradigm and **encodes reasoning operations**

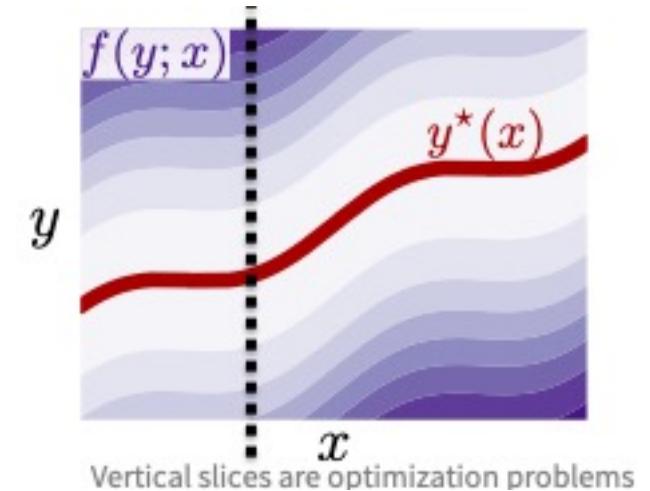
- Finds the **best way to interact** with a **representation of the world**

Focus: parametric optimization problems that are **repeatedly solved**



$$y^*(x) \in \operatorname{argmin}_{y \in \mathcal{C}(x)} f(y; x)$$

optimal solution objective context (or parameterization)
 | |
 | |
 | |
optimization variable constraints



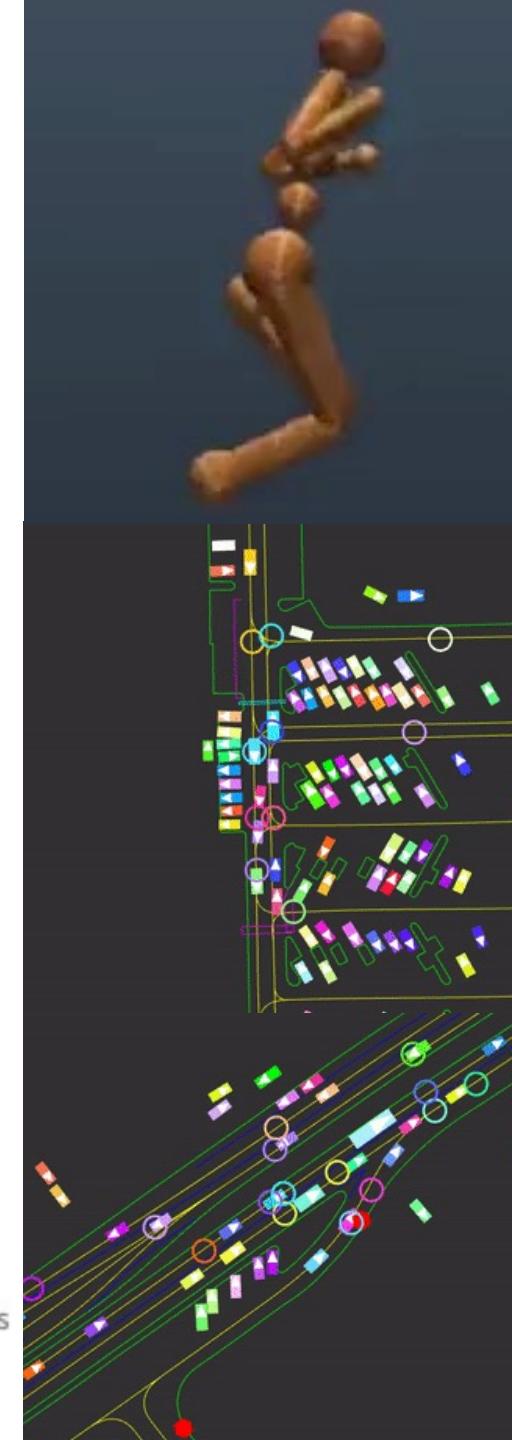
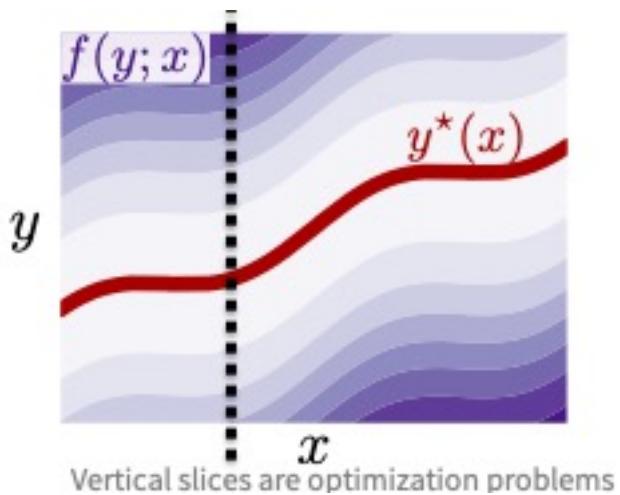
Breakthroughs enabled by optimization include

1. **controlling systems** (robotic, autonomous, mechanical, and multi-agent)
2. **making operational decisions** based on future predictions
3. efficiently **transporting** or **matching** resources, information, and measures for **optimal transport** and **generative modeling**
4. **allocating** budgets and portfolios,
5. **designing** materials, molecules, and other structures,
6. **solving inverse problems** (infers underlying hidden costs, incentives, geometries, terrains, and other structures)
7. **parameter learning** of predictive and statistical models

optimal solution objective context (or parameterization)

$$y^*(x) \in \operatorname{argmin}_{y \in \mathcal{C}(x)} f(y; x)$$

optimization variable constraints



When optimization fails, machine learning helps

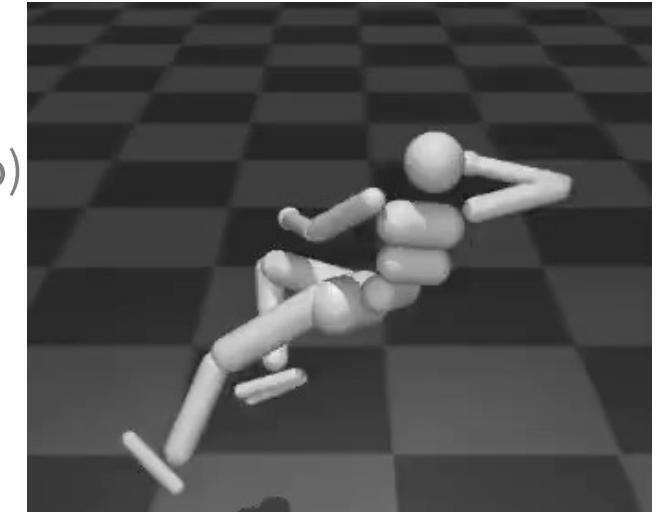
$$y^*(x) \in \operatorname{argmin}_{y \in \mathcal{C}(x)} f(y; x)$$

Optimization starts failing us when:

1. Analytically **encoding every detail** into f and \mathcal{C} is impossible
 - May be **unknown, mis-specified, or inaccurate**
 - Especially difficult with **high-dimensional data** (text, images, video)
2. **Solving** for $y^*(x)$ is **intractable** and **difficult to compute**

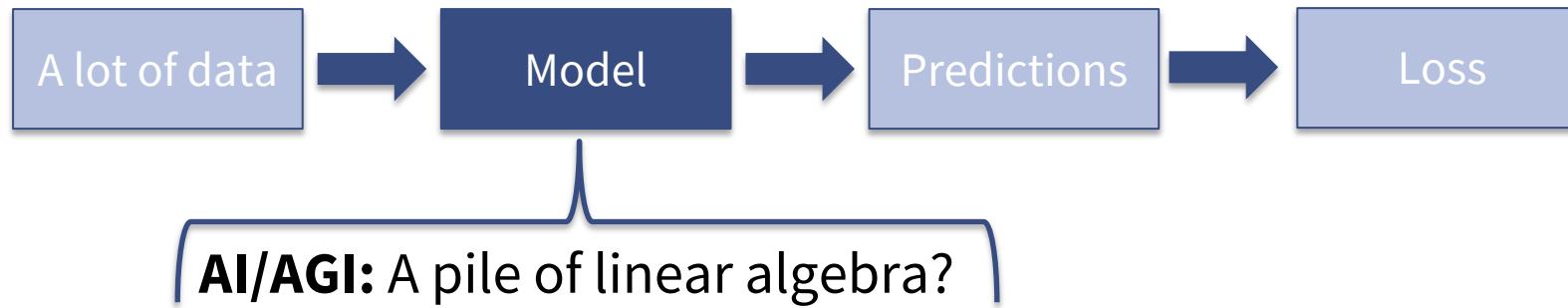
Machine learning systems are **part of the solution** and **excel** at:

1. **Modeling systems** we can extract a lot of data from
2. Extracting **meaningful latent representations**
3. Making **fast and accurate** predictions



When optimization fails, machine learning helps

When machine learning fails, optimization helps

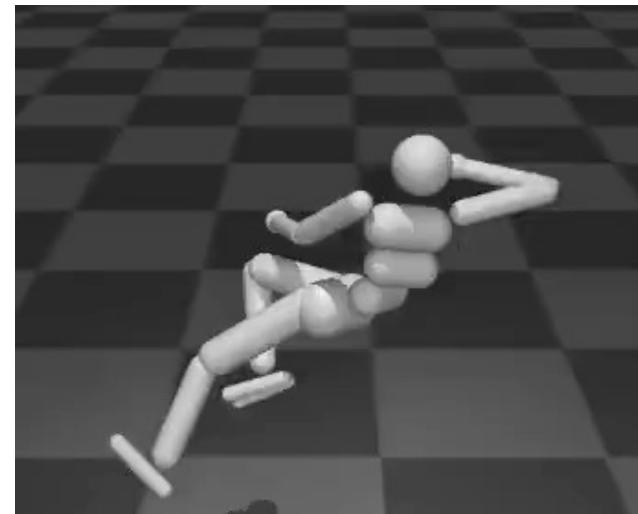


Is a **black-box of linear algebra** the **solution to AI/AGI?** *Maybe...*

Without enough data, learning systems struggle to:

1. **Encode hard constraints** and **rules** of the world crucial for safety
2. Perform **reasoning operations**
3. Isolate **interpretable** and **modular subcomponents**
4. **Generalize** to new settings beyond the training distribution

Optimization often **excels in these** and **can be brought into the model/loss**



Integrating optimization and machine learning

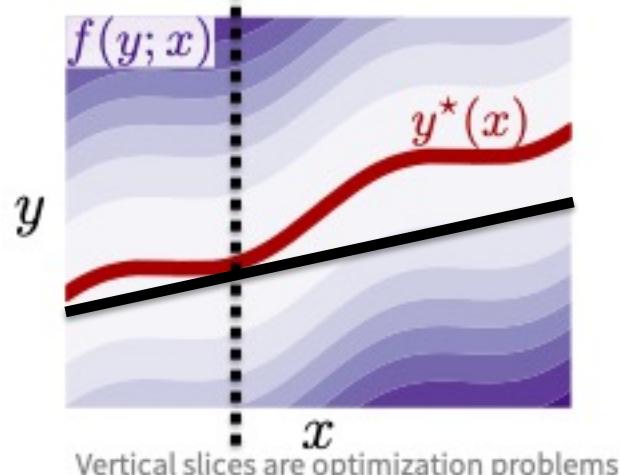
Key: view **optimization as a function** from the context x to the solution $y^*(x) \in \operatorname{argmin}_{y \in \mathcal{C}(x)} f(y; x)$

Differentiable optimization — $\frac{\partial}{\partial x} y^*(x)$

Optimization is a function that can be **differentiated**

Tells how **perturbing the context changes the solution**

Enables **optimization as a layer** or loss for ML

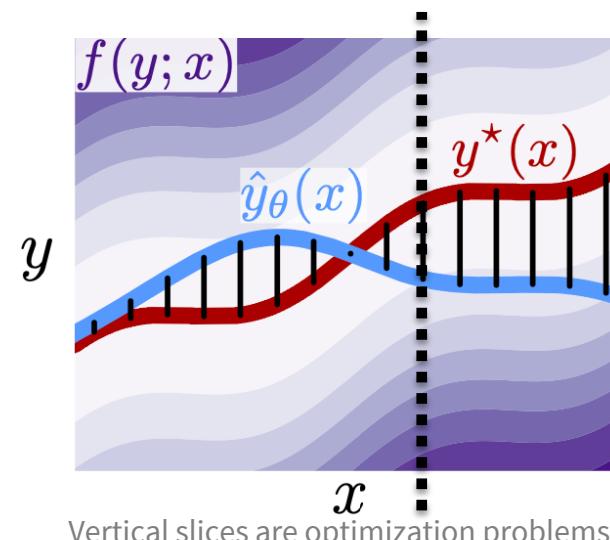


Amortized optimization — $\hat{y}_\theta(x) \approx y^*(x)$

Optimization is a function that can be **approximated**

Approximate the solution map with \hat{y}_θ

Power of optimization, fast speed of an ML model



This talk

Differentiable optimization — $\frac{\partial}{\partial x} y^*(x)$

Foundations

Differentiable convex optimization layers

Differentiable control

Optimization-based task losses

Amortized optimization — $\hat{y}_\theta(x) \approx y^*(x)$

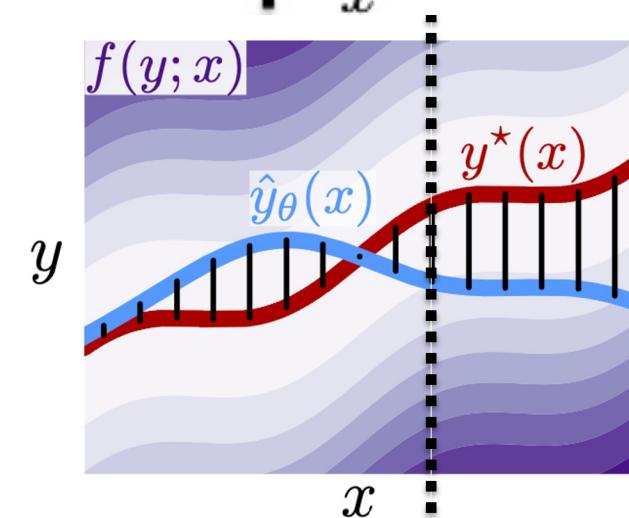
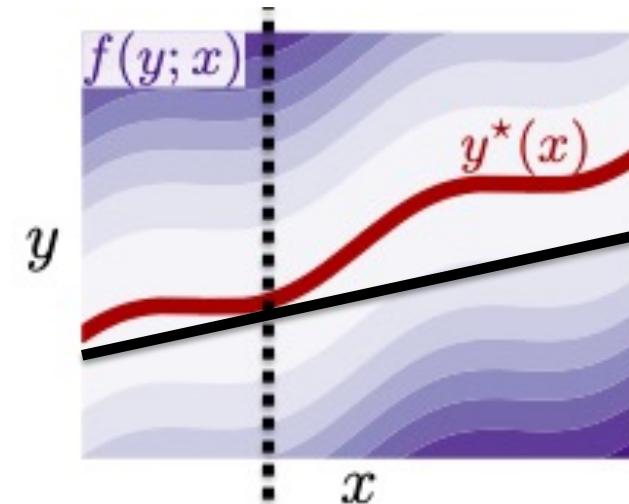
Foundations

RL as amortized optimization

Amortization via learning latent subspaces

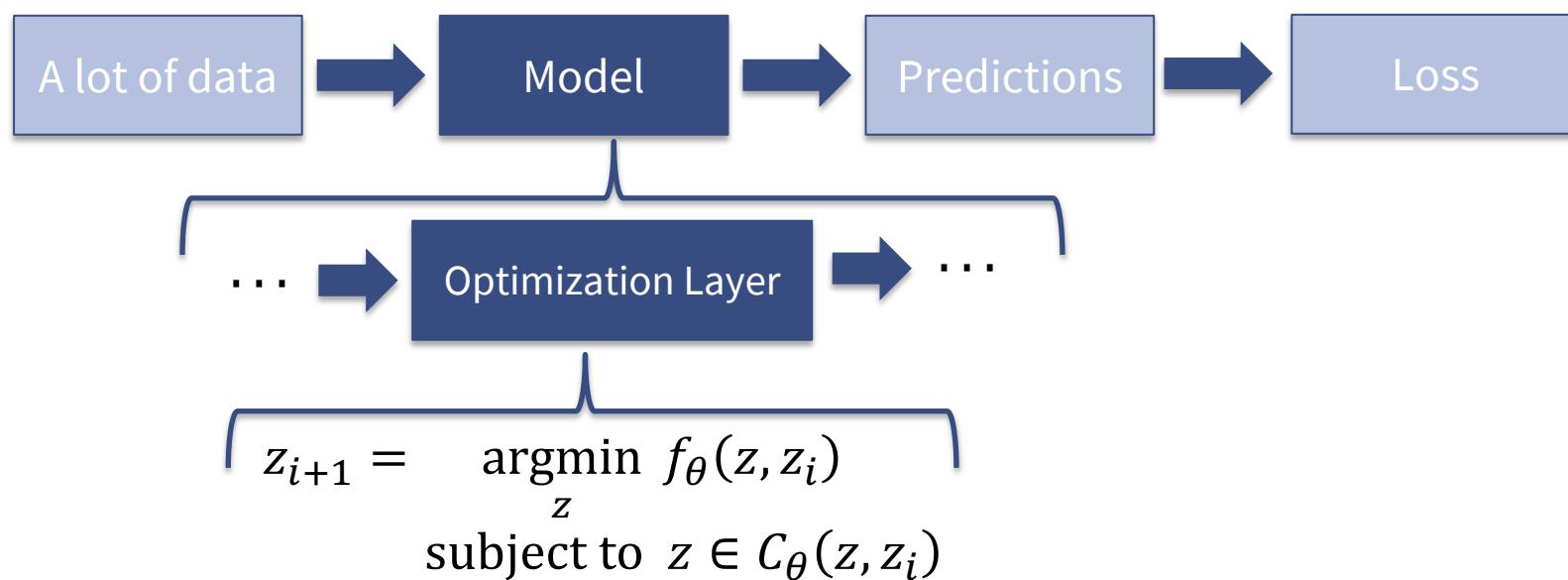
Meta Optimal Transport

Amortizing convex conjugates



Differentiable optimization layers

- **Definition.** A **differentiable optimization layer** for a machine learning model internally solves an optimization problem and is learned with backpropagation
- **Why optimization?** Many **useful operations** can be expressed with optimization (sorting, permutations, top-k, control, planning, power generation, and others from earlier)
- **Why differentiable optimization?** End-to-end learn **through the optimization procedure** along with the other parts of the model using derivatives



Differentiable convex optimization layers

- The **argmin** of a **convex** optimization problem is **non-convex** and **expressive**
- Most models **already** use differentiable convex optimization layers!*

*Useful as an introduction for differentiable convex optimization, not for improving these models

- Standard non-linearities to be seen as **solutions** to convex optimization problem
(ReLU, sigmoid, softmax)

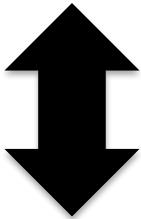
$$y^*(x) = \underset{y}{\operatorname{argmin}} f(y; x) \text{ subject to } y \in C(x)$$

The ReLU is a convex optimization layer

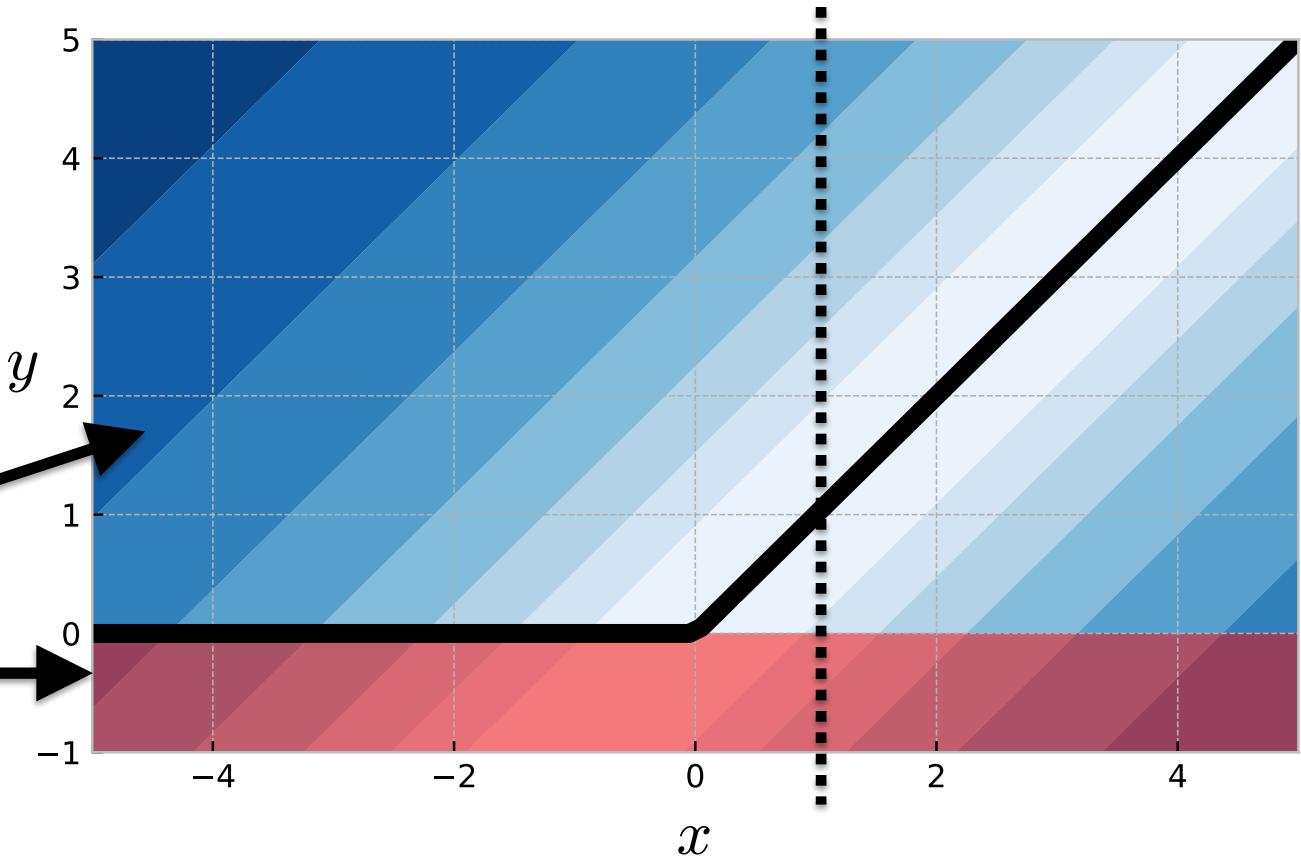
Section 2 of Differentiable optimization-based modeling for machine learning. Amos, PhD Thesis 2019

Proof: Comes from first-order optimality

$$\text{ReLU}(x) = \max\{0, x\}$$



$$\begin{aligned} \text{ReLU}(x) &= \underset{y}{\operatorname{argmin}} \quad \|y - x\|_2^2 \\ \text{s.t. } y &\geq 0 \end{aligned}$$

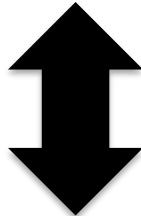


The sigmoid is a convex optimization layer

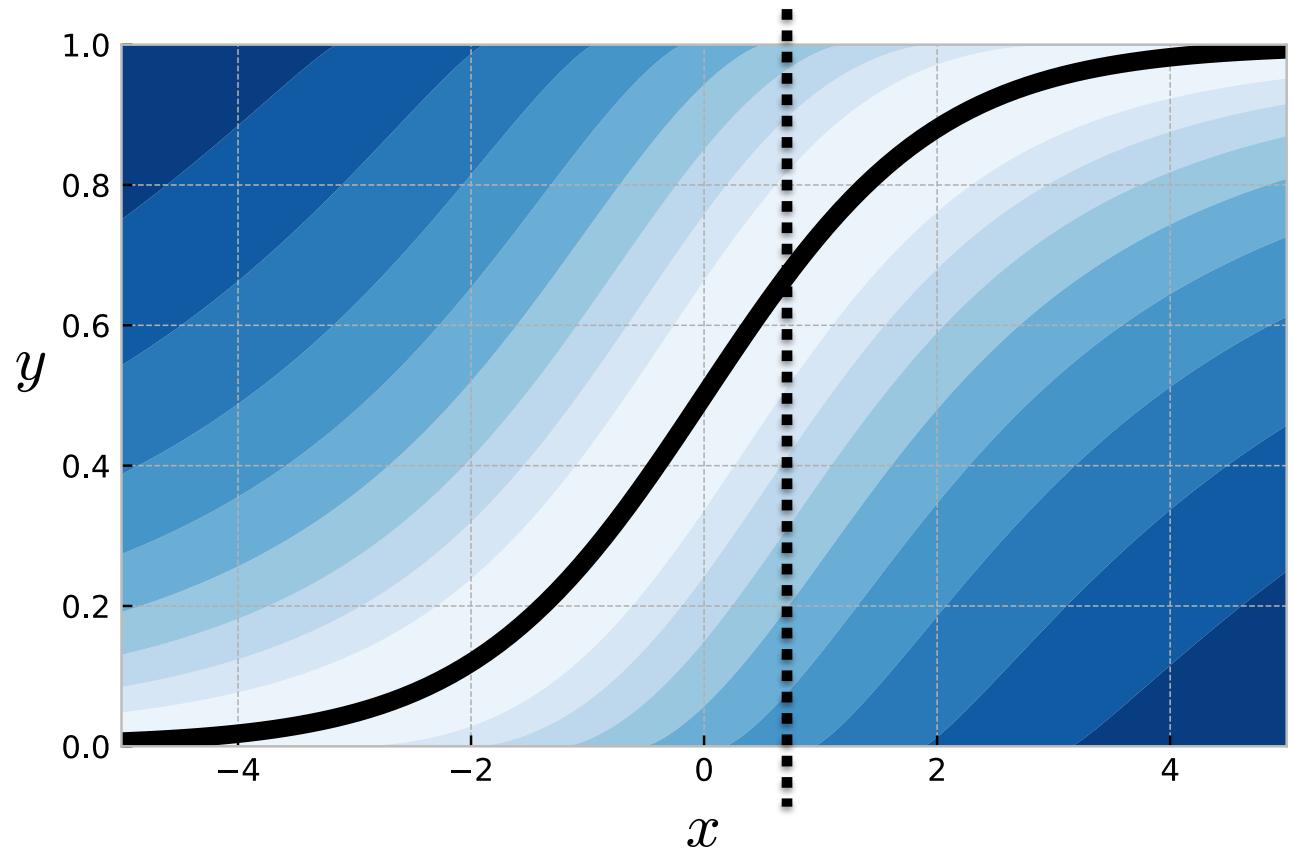
Section 2 of Differentiable optimization-based modeling for machine learning. Amos, PhD Thesis 2019

Proof: Comes from first-order optimality

$$\sigma(x) = \frac{1}{1 + \exp \{-x\}}$$



$$\begin{aligned} \sigma(x) = \operatorname{argmin}_y & -y^T x - H_b(y) \\ \text{s.t. } & 0 \leq y \leq 1 \end{aligned}$$



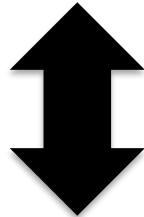
$H_b(y) := -\sum_i(y_i \log y_i + (1 - y_i) \log(1 - y_i))$ is the binary cross-entropy function

The softargmax is a convex optimization layer

Section 2 of *Differentiable optimization-based modeling for machine learning*. Amos, PhD Thesis 2019

Proof: Comes from first-order optimality

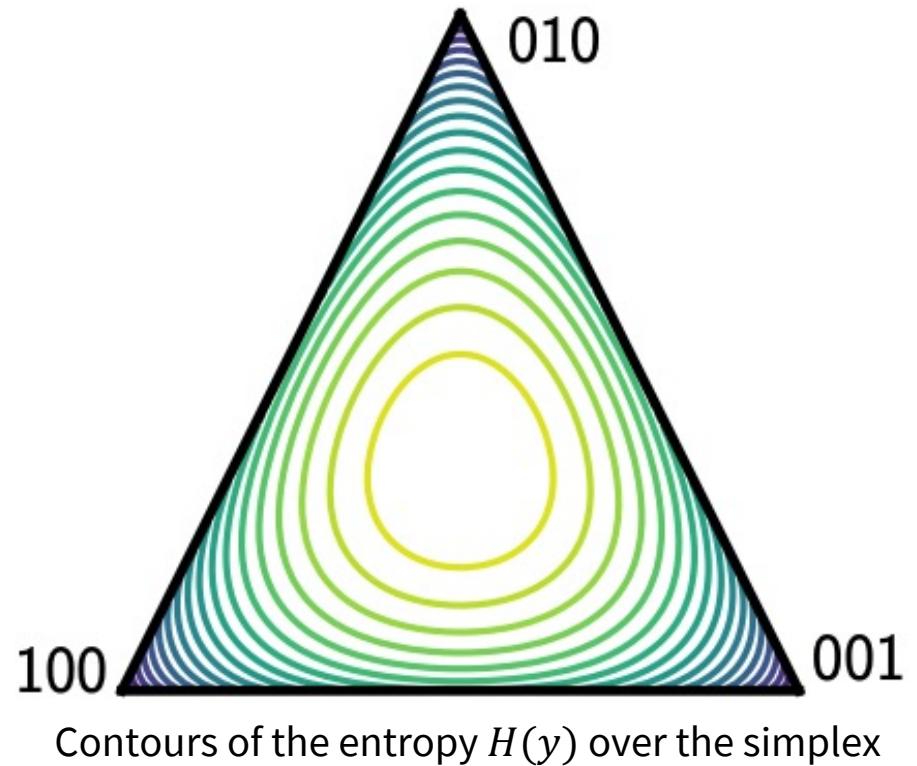
$$\pi_{\Delta}(x) = \frac{\exp x}{\sum_i \exp x_i}$$



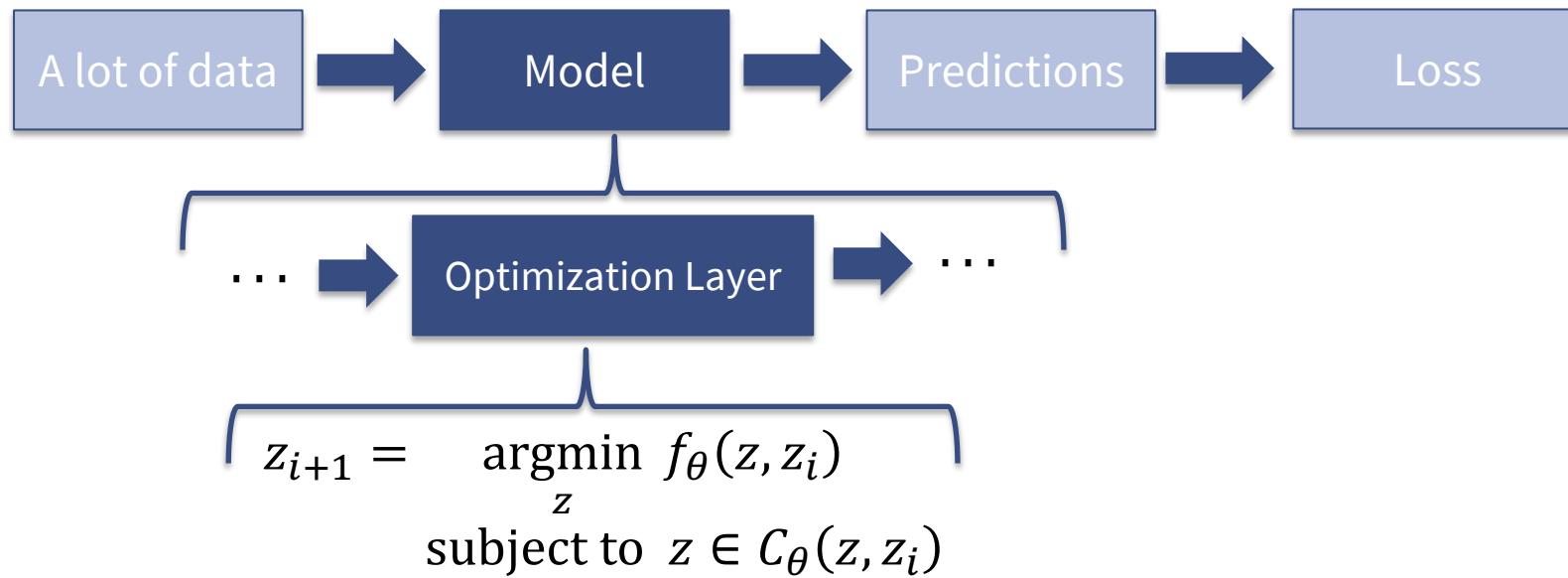
$$\pi_{\Delta}(x) = \operatorname{argmin}_y -y^T x - H(y)$$

$$\begin{aligned} \text{s.t. } & 0 \leq y \leq 1 \\ & 1^T y = 1 \end{aligned}$$

$H(y) := -\sum_i y_i \log y_i$ is the entropy function



How can we generalize this?



Derivatives and backpropagation

For learning, we **differentiate** or backpropagate through these layers — **differentiable optimization**

Easy if the optimization problem has an **explicit, closed-form solution** (often standard differentiation)

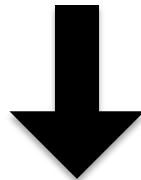
Otherwise, need to use **implicit differentiation**, which is also used for sensitivity analysis

Differentiable convex quadratic programs

OptNet: Differentiable Optimization as a Layer in Neural Networks. Amos and Kolter, ICML 2017.

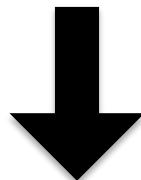
$$x^* = \underset{x}{\operatorname{argmin}} \frac{1}{2} x^\top Q x + p^\top x$$

subject to $Ax = b$ $Gx \leq h$



KKT Optimality

Find z^* s.t. $\mathcal{R}(z^*, \theta) = 0$ where $z^* = [x^*, \dots]$ and $\theta = \{Q, p, A, b, G, h\}$



Implicitly differentiating \mathcal{R} gives $D_\theta(z^*) = -\left(D_z \mathcal{R}(z^*)\right)^{-1} D_\theta \mathcal{R}(z^*)$

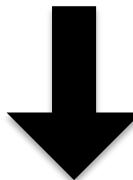
Differentiable convex conic programs

Section 7 of Differentiable optimization-based modeling for machine learning. Amos, PhD Thesis 2019

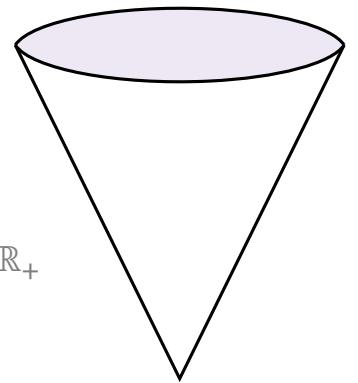
Differentiating through a cone program. Agrawal et al., 2019

Differentiable convex optimization layers. Agrawal*, Amos*, Barratt*, Boyd*, Diamond*, Kolter*, NeurIPS 2019.

$$x^* = \underset{x}{\operatorname{argmin}} \quad c^T x \\ \text{subject to } b - Ax \in \mathcal{K}$$

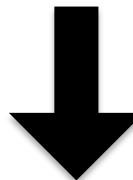


- Zero: $\{0\}$
- Free: \mathbb{R}^n
- Non-negative: \mathbb{R}_+^n
- Second-order (Lorentz): $\{(t, x) \in \mathbb{R}_+ \times \mathbb{R}^n \mid \|x\|_2 \leq t\}$
- Semidefinite: \mathbb{S}_+^n
- Exponential: $\{(x, y, z) \in \mathbb{R}^3 \mid ye^{x/y} \leq z, y > 0\} \cup \mathbb{R}_- \times \{0\} \times \mathbb{R}_+$
- Cartesian Products: $\mathcal{K} = \mathcal{K}_1 \times \dots \times \mathcal{K}_p$



Conic Optimality

Find z^* s.t. $\mathcal{R}(z^*, \theta) = 0$ where $z^* = [x^*, \dots]$ and $\theta = \{A, b, c\}$



Implicitly differentiating \mathcal{R} gives $D_\theta(z^*) = -(D_z \mathcal{R}(z^*))^{-1} D_\theta \mathcal{R}(z^*)$

From the softmax to soft and differentiable top-k

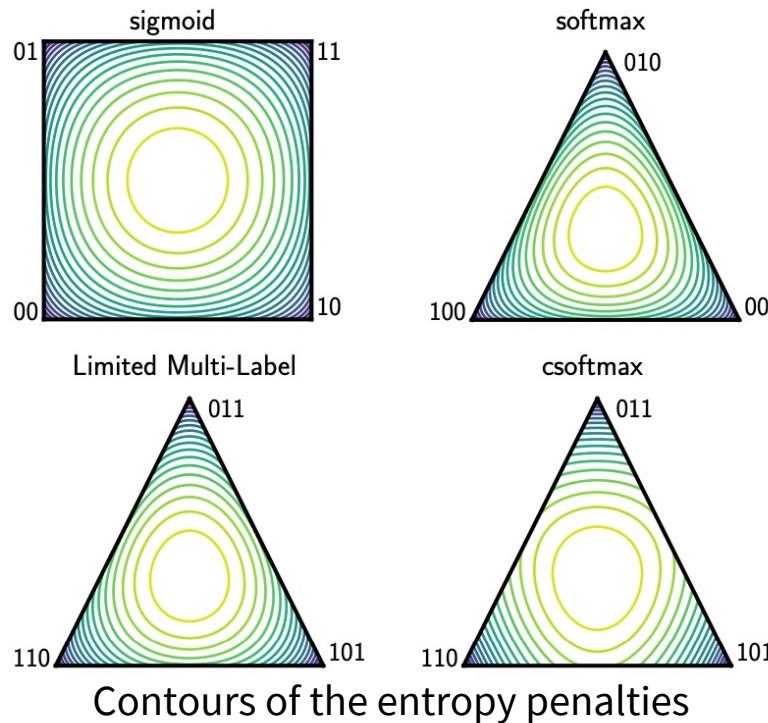
Limited multi-label projection layer. Amos et al., 2019.

$$y^* = \underset{y}{\operatorname{argmin}} -y^\top x - H(y)$$

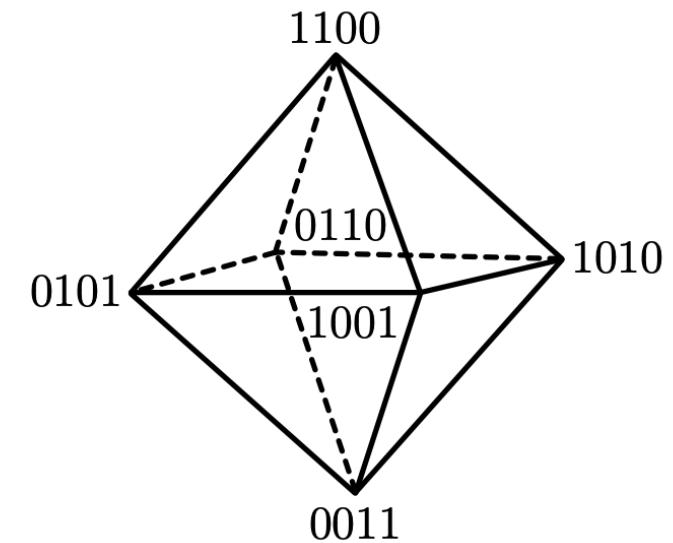
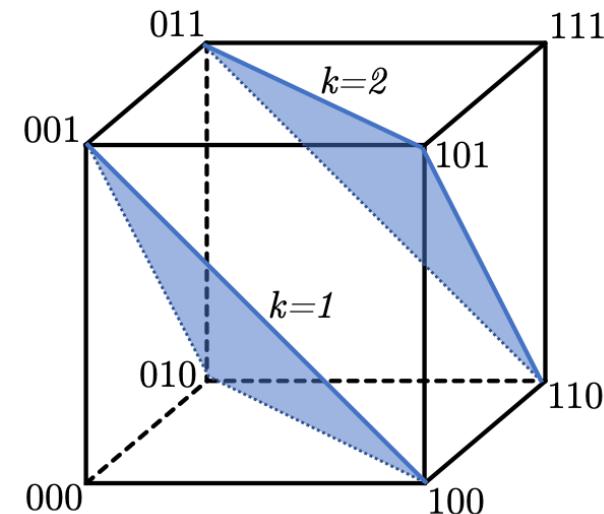
subject to $0 \leq y \leq 1$

$$1^\top y = 1$$

$$\begin{aligned} y^* = \operatorname*{argmin}_y \quad & -y^\top x - H_b(y) \\ \text{subject to} \quad & 0 \leq y \leq 1 \\ & 1^\top y = k \end{aligned}$$



$H_b(y) \coloneqq -\sum_i(y_i \log y_i + (1 - y_i) \log(1 - y_i))$ is the binary cross-entropy function



Differentiable permutations, sorting and SVMs

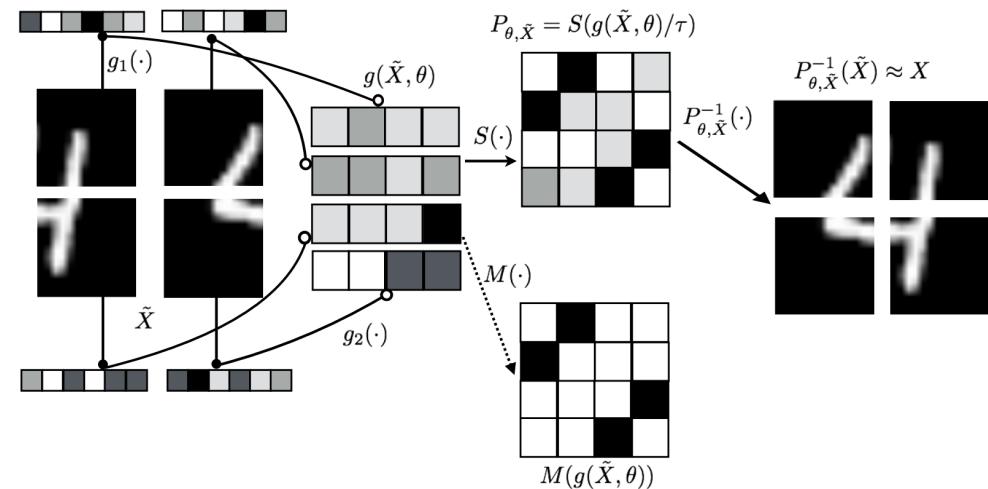
Differentiable permutations and sorting (Gumbel-Sinkhorn)

Projection onto the **Birkhoff polytope** \mathcal{B}_N :

$$S\left(\frac{X}{\tau}\right) = \underset{P \in \mathcal{B}_N}{\operatorname{argmax}} \langle P, X \rangle_F + \tau H(P)$$

$$\mathcal{B}_N = \{X: X \geq 0, \sum_i X_{ij} = \sum_j X_{ij} = 1\}$$

Learning latent permutations with Gumbel-Sinkhorn networks. Mena et al., ICLR 2018.

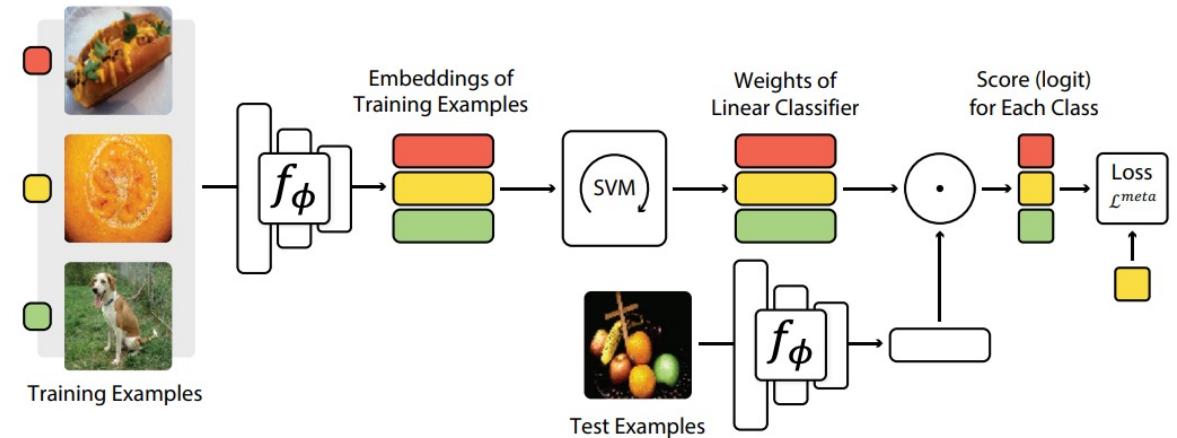


Differentiable SVMs (MetaOptNet)

Differentiate the decision boundary w.r.t. the dataset

$$w^*(\mathcal{D}) = \underset{w}{\operatorname{argmin}} \|w\|^2 + C \sum_i \max\{0, 1 - y_i f(x_i)\}$$

Meta-learning with differentiable convex optimization. Lee et al., CVPR 2019.



Sensitivity and perturbation analysis

Differentiable convex optimization layers. Agrawal*, Amos*, Barratt*, Boyd*, Diamond*, Kolter*, NeurIPS 2019.

Adjoint derivatives for optimization problems have been studied for decades
We have focused on uses for learning, but also widely used for **sensitivity analysis**

Logistic regression example

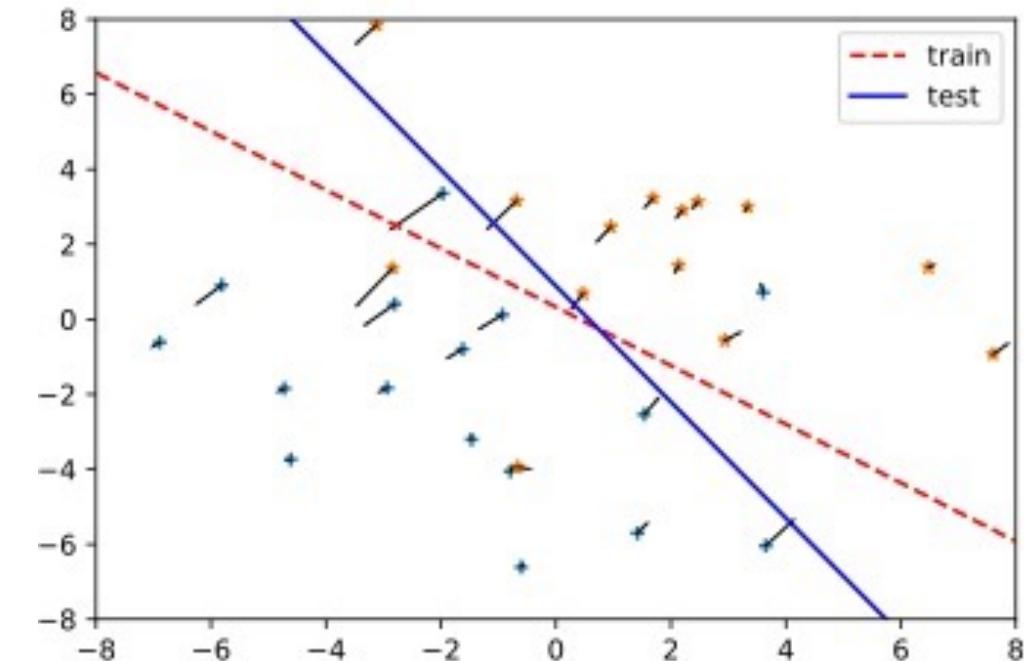
Find optimal decision boundary:

$$\theta^*(\mathcal{D}) \in \operatorname{argmax}_{\theta} \sum_i \log p_{\theta}(y_i | x_i)$$

$\theta^*(\mathcal{D})$ is **just a function** of the data

Derivatives $\frac{\partial \theta^*}{\partial x_i}$ provide **sensitivity** to the data points

- *Show how much the data impacts the decision boundary*



Optimal control drives systems

Setting: deterministic, discrete-time system with a **continuous state-action space**

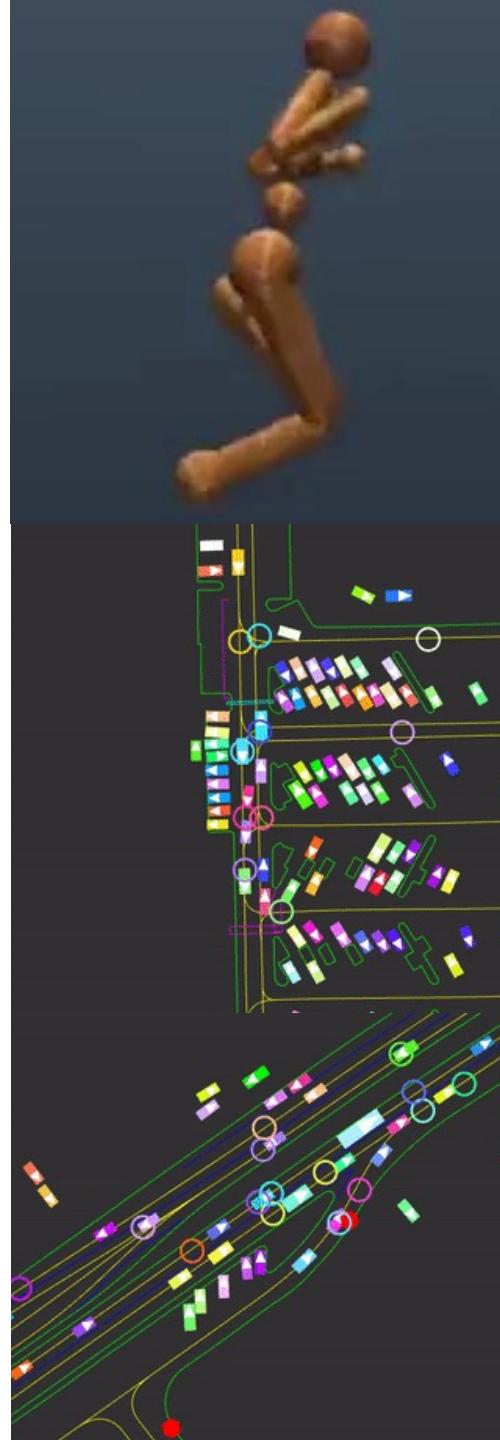
$$x_{1:T}^*, u_{1:T}^* \in \operatorname{argmin}_{x_{1:T}, u_{1:T}} \sum_t \begin{array}{|c|} \text{cost} \\ \hline C(x_t, u_t) \end{array} \text{ s.t. } \begin{array}{|c|} \text{initial state} \\ \hline x_1 = x_{\text{init}} \end{array} \quad \begin{array}{|c|} \text{dynamics} \\ \hline x_{t+1} = f(x_t, u_t) \end{array} \quad \begin{array}{|c|} \text{constraints} \\ \hline u_t \in \mathcal{U} \end{array}$$

Full notation: $u_{1:T}^*(x_{\text{init}}, f_\theta)$

Widely deployed over the past century in aviation, robotics, vehicles, HVAC
Often for a **Markov decision process** but doesn't have to be

The **real-world is non-convex**, so are our controllers
Convex in some cases, e.g., with quadratic cost/linear dynamics (LQR)

No learning necessary if we know the system — just pure optimization



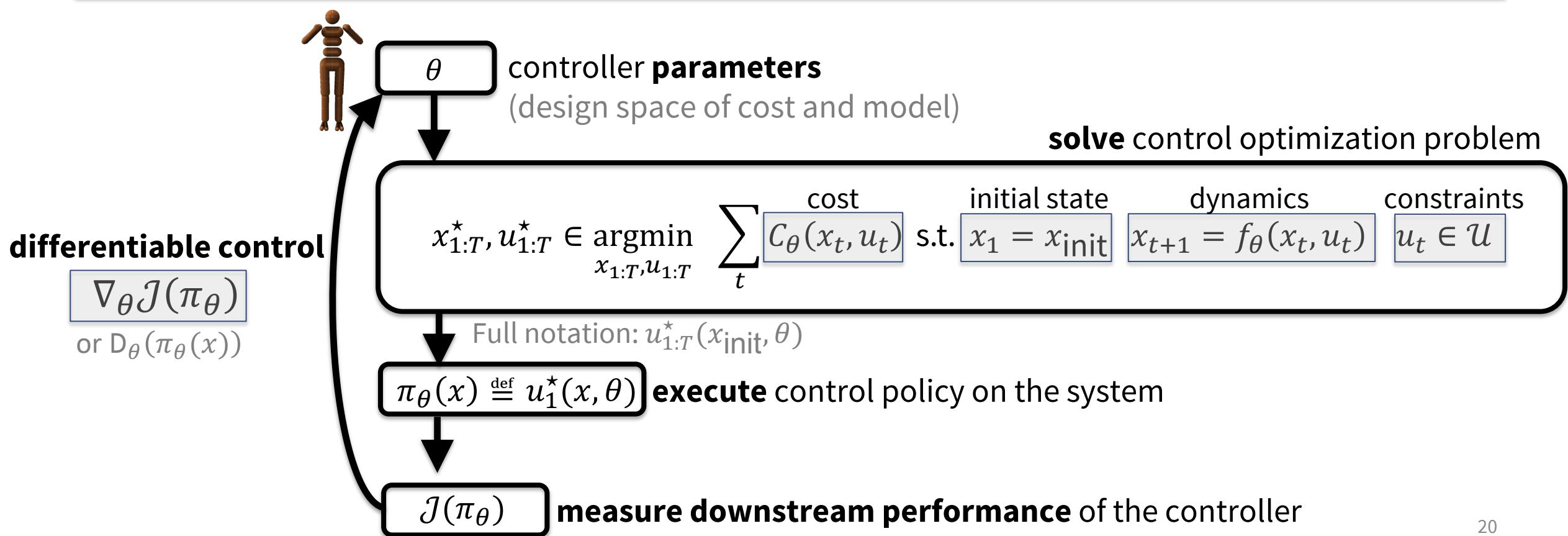
Differentiable control

Differentiable MPC for end-to-end planning and control. Amos et al., NeurIPS 2018.

The differentiable cross-entropy method. Amos and Yarats, ICML 2020.

We can often measure the **downstream performance** induced by the controller

Idea: optimize (i.e., tune/learn) the parameters for a **downstream performance metric** by **differentiating through the control optimization problem**



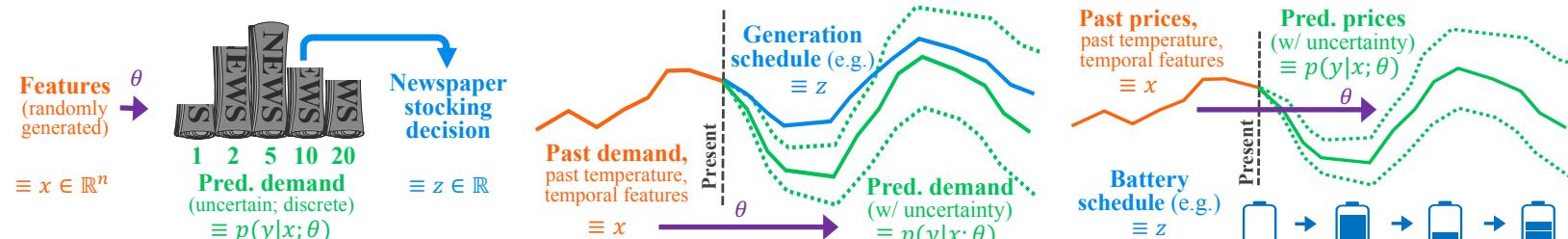
Optimization as a loss for downstream tasks

Task-based end-to-end model learning in stochastic optimization. Donti, Amos, and Kolter, NeurIPS 2017.

A model's predictions may be used in a **downstream optimization problem** (e.g., scheduling)

Challenge: models are **inaccurate** and **not aware of how errors impact the downstream problem**

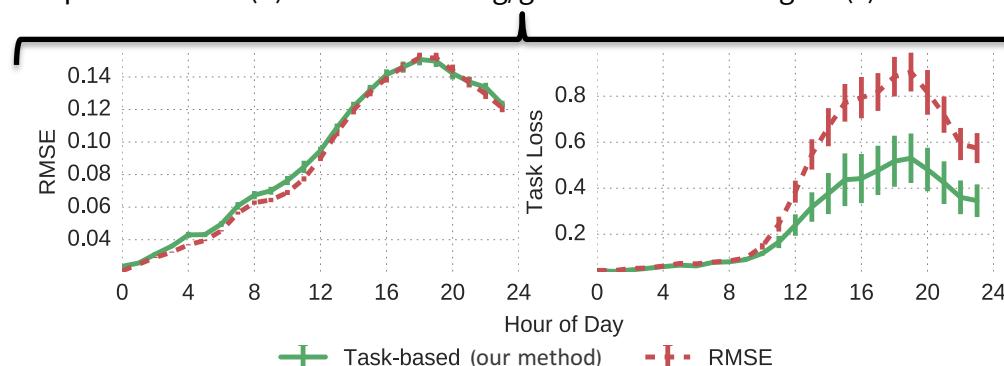
Idea: use the **downstream optimizer as a loss** for the model using differentiable optimization



(a) Inventory stock problem

(b) Load forecasting/generator scheduling

(c) Price forecasting/battery arbitrage



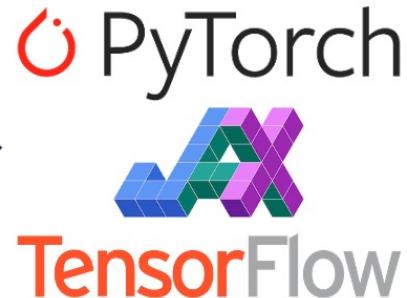
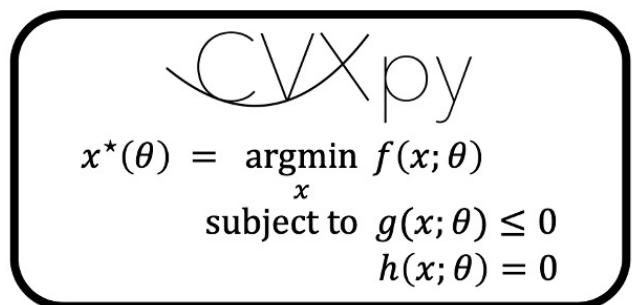
$$\underset{z \in \mathbb{R}^{24}}{\text{minimize}} \mathbf{E}_y \left[\sum_{i=1}^{24} \gamma_s [y_i - z_i]_+ + \gamma_e [z_i - y_i]_+ + \frac{1}{2} \|z_i - y_i\|^2 \right] \quad \text{subject to } |z_i - z_{i+1}| \leq c_{\text{ramp}}.$$

Implementing differentiable optimization

Manually: unroll or implicitly differentiate the optimality conditions (such as the KKT system)
My original differentiable QP implementation was **~1k lines of code** to make efficient

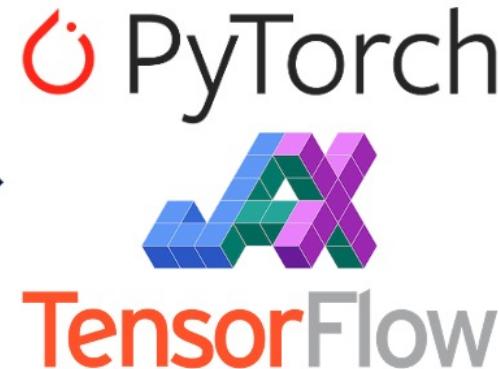
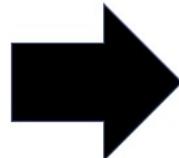
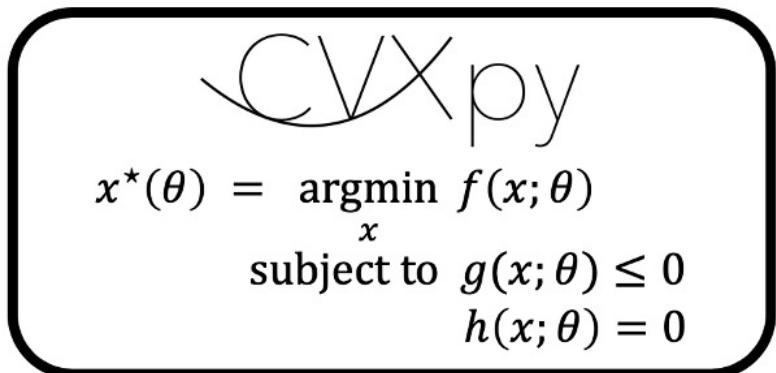
Use a software package to prototype such as:

- **cvxpylayers**: differentiable convex optimization in CVXPY [Agrawal*, Amos*, et al., NeurIPS 2019]
Implements a differentiable QP in ~10 lines of code
- **Theseus**: differentiable non-linear least squares with a focus on robotics [Pineda et al., NeurIPS 2022]
- **higher**: differentiable unrolled optimizers [Grefenstette et al., 2019]



Differentiable CVXPY

Differentiable convex optimization layers. Agrawal*, Amos*, Barratt*, Boyd*, Diamond*, Kolter*, NeurIPS 2019.



(Officially part of CVXPY!)

$$z_{i+1} = \underset{z}{\operatorname{argmin}} \frac{1}{2} z^\top Q(z_i) z + q(z_i)^\top z$$

subject to

$$A(z_i)z = b(z_i)$$
$$G(z_i)z \leq h(z_i)$$

Parameters/Submodules : Q, q, A, b, G, h

Before: 1k lines of code, **now:**

```
import cvxpy as cp
from cvxpyth import CvxpyLayer
obj = cp.Minimize(0.5*cp.quad_form(x, Q) + p.T * x)
cons = [A*x == b, G*x <= h]
prob = cp.Problem(obj, cons)
layer = CvxpyLayer(prob, params=[Q, p, A, b, G, h], out=[x])
```

This talk

Differentiable optimization — $\frac{\partial}{\partial x} y^*(x)$

Foundations

Differentiable convex optimization layers

Differentiable control

Optimization-based task losses

Amortized optimization — $\hat{y}_\theta(x) \approx y^*(x)$

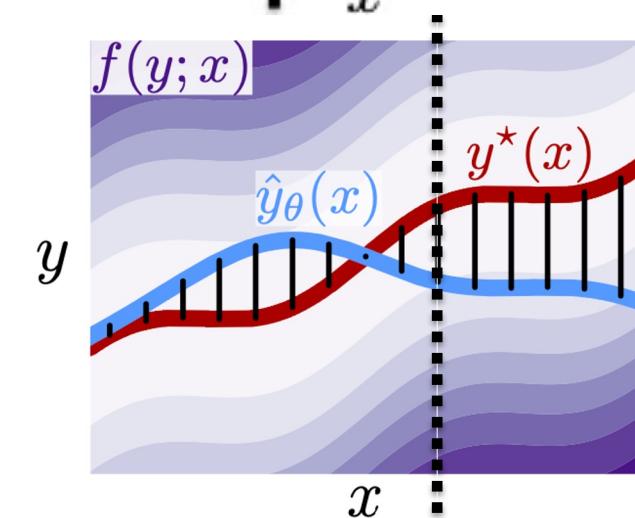
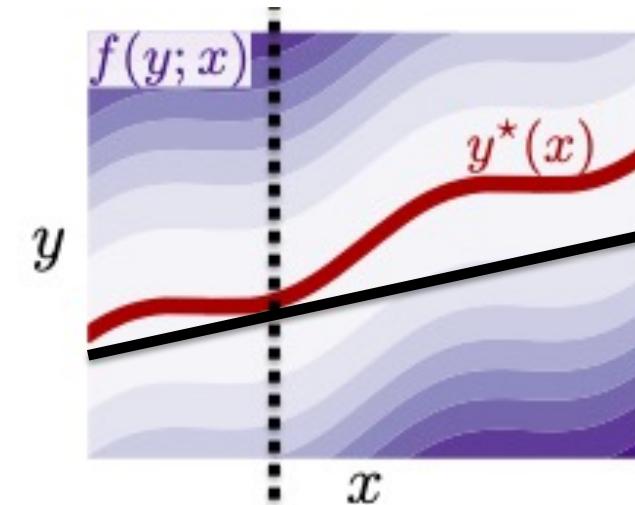
Foundations

RL as amortized optimization

Amortization via learning latent subspaces

Meta Optimal Transport

Amortizing convex conjugates



Amortized optimization and learning to optimize

Tutorial on amortized optimization for learning to optimize over continuous domains. Amos, Foundations and Trends in Machine Learning (to appear)

- **Difficulty:** optimization is computationally expensive
 - Sometimes **solving once** may be difficult
 - Exasperated when **repeatedly** solving, e.g., as part of a differentiable optimization layer
- **Insight:** optimal solutions share structure and don't live in isolation
- **Solution:** amortized optimization (or learning to optimize)
 - Use machine learning to uncover the shared structure
 - Create learning-augmented versions of classical optimization solvers
 - Far surpasses average or worst-case convergence rates

Amortized variational inference and VAEs

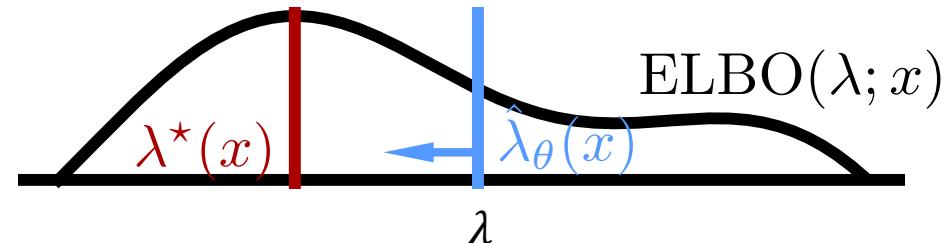
Origin of the term “amortization”

Given a VAE model $p(x) = \log \int_z p(x|z)p(z)$, **encoding** can be done by solving the optimization problem

$$\lambda^*(x) = \underset{\lambda}{\operatorname{argmax}} \text{ELBO}(\lambda; x) \quad \text{where} \quad \text{ELBO}(\lambda; x) := \mathbb{E}_{q(z;\lambda)}[\log p(x|z)] - D_{\text{KL}}(q(x;\lambda)||p(z)).$$

VAEs amortize this (along with other amortized variational inference methods) and solve:

$$\underset{\theta}{\operatorname{argmax}} \mathbb{E}_{x \sim p(x)} \text{ELBO}(\hat{\lambda}_\theta; x)$$



Amortization, optimal control, and RL

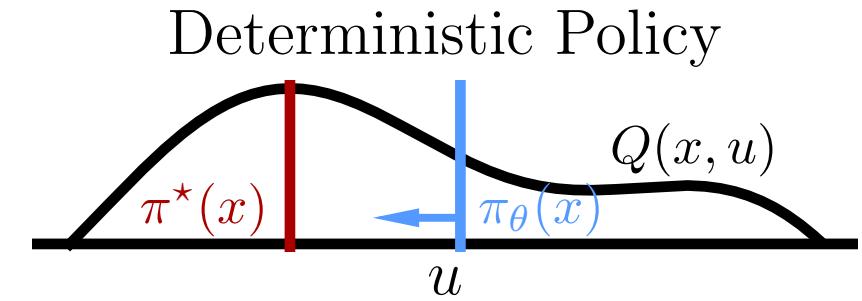
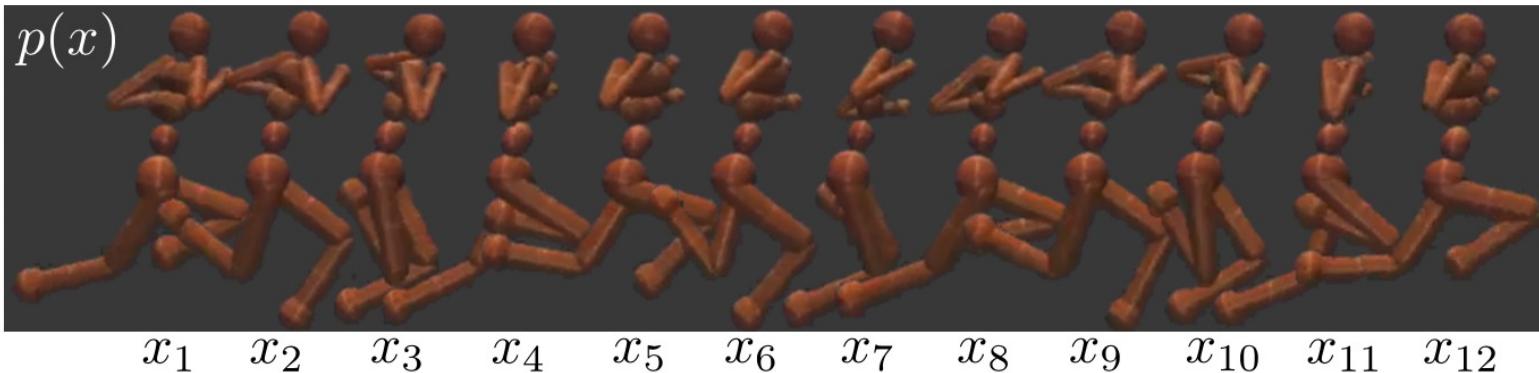
Tutorial on amortized optimization for learning to optimize over continuous domains. Amos, Foundations and Trends in Machine Learning (to appear)
On the model-based stochastic value gradient for continuous reinforcement learning. Amos et al., L4DC 2021.

Given a value estimate $Q(x, u)$, **optimal control** finds the best action u for a given state x by solving

$$\pi^*(x) := \underset{u}{\operatorname{argmax}} Q(x, u)$$

Almost every reinforcement learning method amortizes this optimization problem!

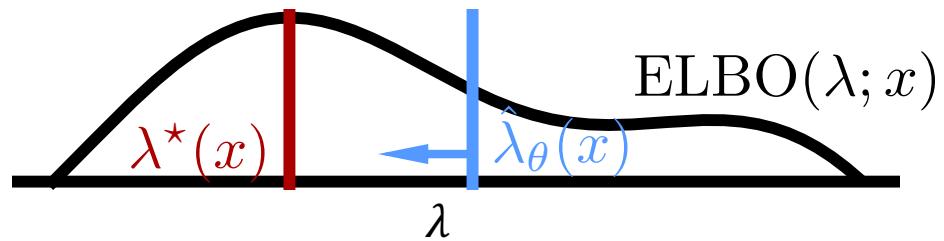
$$\underset{\theta}{\operatorname{argmax}} \mathbb{E}_{p(x)} Q(x, \pi_{\theta}(x))$$



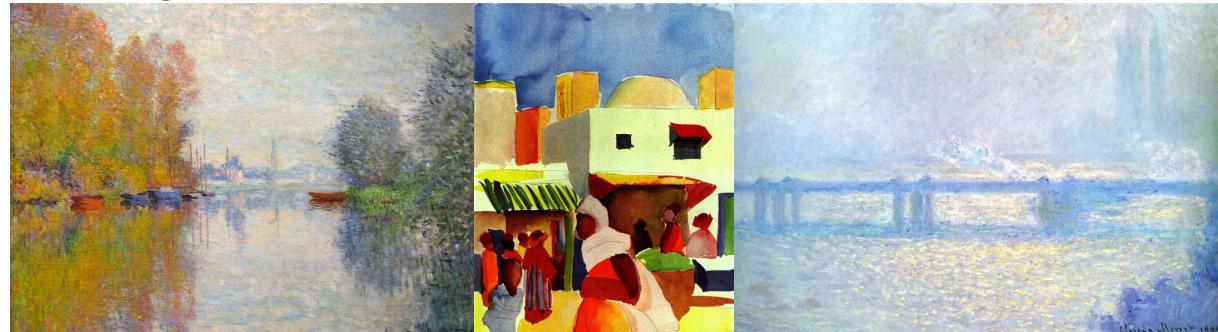
Amortization similarity between VAEs and RL

VAE posterior amortization

$$\operatorname{argmax}_{\theta} \mathbb{E}_{p(x)} \text{ELBO}(\hat{\lambda}_{\theta}(x); x)$$



x : images from dataset



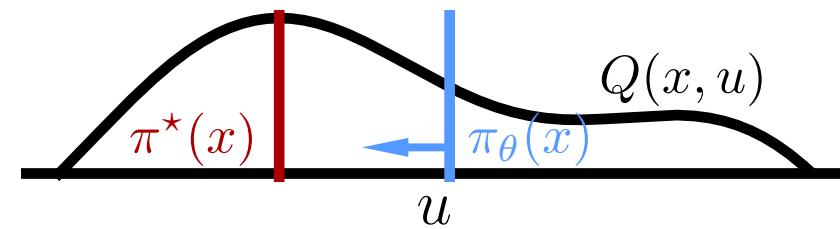
x_1

x_2

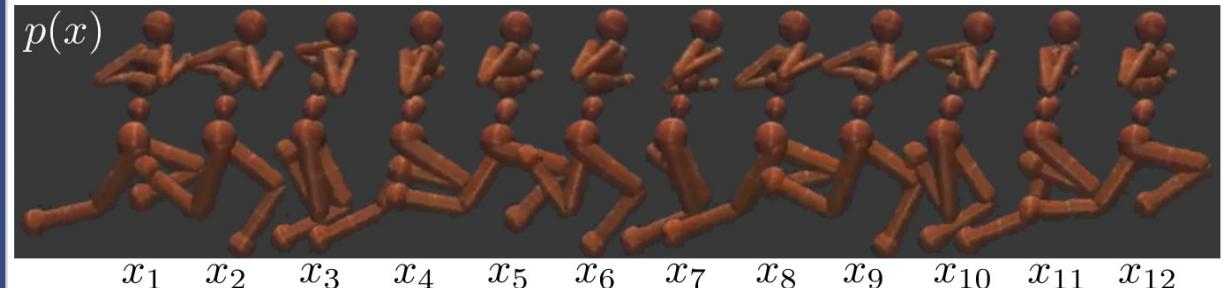
x_3

Value gradient amortization in RL

$$\operatorname{argmax}_{\theta} \mathbb{E}_{p(x)} Q(x, \pi_{\theta}(x))$$



x : states from system



$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \quad x_7 \quad x_8 \quad x_9 \quad x_{10} \quad x_{11} \quad x_{12}$

Generalizing amortization beyond these

Foundations and Trends® in Machine Learning

Tutorial on amortized optimization for learning to optimize over continuous domains

Brandon Amos

Facebook AI Research, Meta

BDA@FB.COM

Abstract

Optimization is a ubiquitous modeling tool and is often deployed in settings which repeatedly solve similar instances of the same problem. Amortized optimization methods use learning to predict the solutions to problems in these settings. This leverages the shared structure between similar problem instances. In this tutorial, we will discuss the key design choices behind amortized optimization, roughly categorizing 1) models into fully-amortized and semi-amortized approaches, and 2) learning methods into regression-based and objective-based. We then view existing applications through these foundations to draw connections between them, including for manifold optimization, variational inference, sparse coding, meta-learning, control, reinforcement learning, convex optimization, and deep equilibrium networks. This framing enables us easily see, for example, that the amortized inference in variational autoencoders is conceptually identical to value gradients in control and reinforcement learning as they both use fully-amortized models with an objective-based loss.

Other examples of amortized optimization

Tutorial on amortized optimization for learning to optimize over continuous domains. Amos, Foundations and Trends in Machine Learning (to appear)

My goal: characterize and connect applications previously developed independently

Variational inference		6.1 VAE SAVAE/IVAE	- ELBO 	variational posterior 	data 	full semi	\mathcal{L}_{obj}
Sparse coding		6.2 PSD LISTA	reconstruction 	sparse code 	data 	full semi	\mathcal{L}_{reg}
Meta-learning		6.3 HyperNets LM MAML Neural Potts	task loss pseudo-likelihood	model parameters 	tasks protein sequences	full semi full	\mathcal{L}_{obj} \mathcal{L}_{RL} \mathcal{L}_{obj} \mathcal{L}_{obj} \mathcal{L}_{obj}
Fixed-points and convex optimization		6.4 NeuralFP HyperAA NeuralSCS HyperDEQ NeuralNMF RLQP	FP residual CP residual DEQ residual NMF residual R_{RLQP}	FP iterates CP iterates DEQ iterates factorizations QP iterates	FP contexts CP parameters DEQ parameters input matrices QP parameters	semi	$\mathcal{L}_{\Sigma}^{\text{obj}}$ $\mathcal{L}_{\text{reg}}^{\Sigma}$ $\mathcal{L}_{\Sigma}^{\text{obj}}$ $\mathcal{L}_{\text{reg}}^{\Sigma}$ $\mathcal{L}_{\Sigma}^{\text{obj}}$ $\mathcal{L}_{\text{RL}}^{\text{obj}}$
Optimal transport		6.5 AmorConj \mathcal{A} -SW Meta OT	c -transform obj max-sliced dist dual OT cost	supp(α) slices Θ optimal couplings	supp(β) mini-batches input measures	full	\mathcal{L}_{obj} \mathcal{L}_{obj} \mathcal{L}_{obj}
Reinforcement learning		6.6 (D)DPG/TD3 PILCO POPLIN DCEM IAPO SVG SAC GPS	- Q -value D_Q or - \mathcal{E}_Q 	controls control dists	state space 	full full semi 	\mathcal{L}_{reg} \mathcal{L}_{obj} \mathcal{L}_{obj} \mathcal{L}_{reg} \mathcal{L}_{reg} \mathcal{L}_{obj} \mathcal{L}_{obj} \mathcal{L}_{obj} \mathcal{L}_{KL}

Amortized optimization definitions

Tutorial on amortized optimization for learning to optimize over continuous domains. Amos, Foundations and Trends in Machine Learning (to appear)

Two key components:

1. **Amortization model** $\hat{y}_\theta(x)$ tries to approximate $y^*(x)$

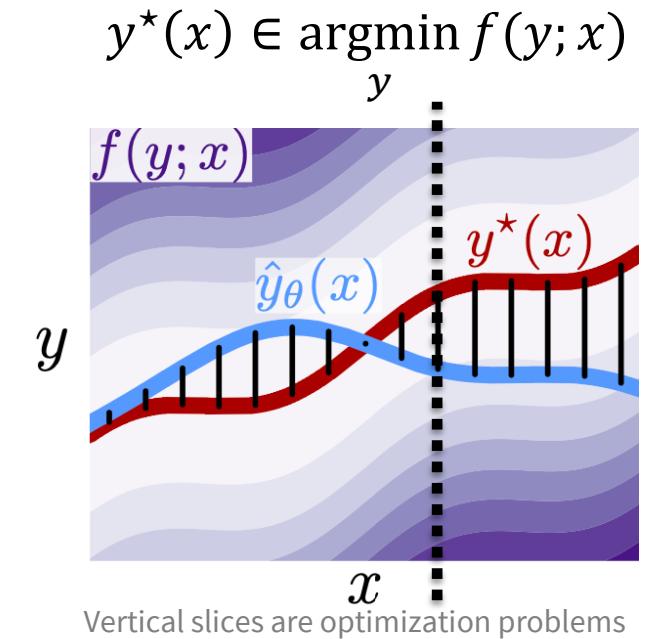
Fully-amortized: A neural network mapping from x to the solution

Semi-amortized: Allowed to query the objective f

2. **Loss** \mathcal{L} measures how well \hat{y} fits y^* and optimized with $\min_\theta \mathcal{L}(\hat{y}_\theta)$

Regression: $\mathcal{L}_{\text{reg}}(\hat{y}_\theta) := \mathbb{E}_{p(x)} \|\hat{y}_\theta(x) - y^*(x)\|_2^2$

Objective-based: $\mathcal{L}_{\text{obj}}(\hat{y}_\theta) := \mathbb{E}_{p(x)} f(\hat{y}_\theta(x); x)$



Amortization is widely deployed (although often not connected as amortization)

Amortized variational inference (VAEs)

Meta-learning (hypernetworks, MAML)

Reinforcement learning (policy learning for actor-critic methods, SAC)

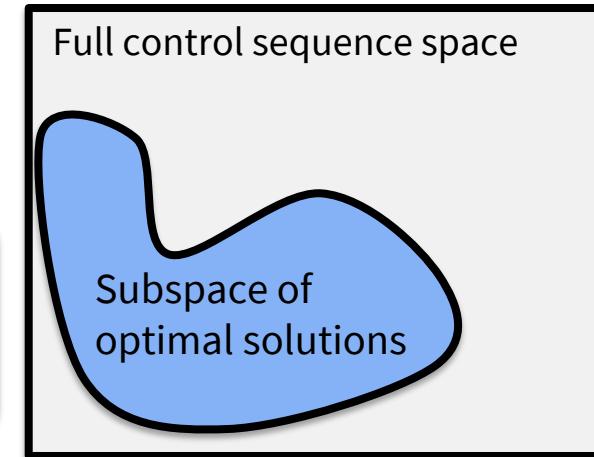
Amortization via learning latent subspaces

The differentiable cross-entropy method. Amos and Yarats, ICML 2020.

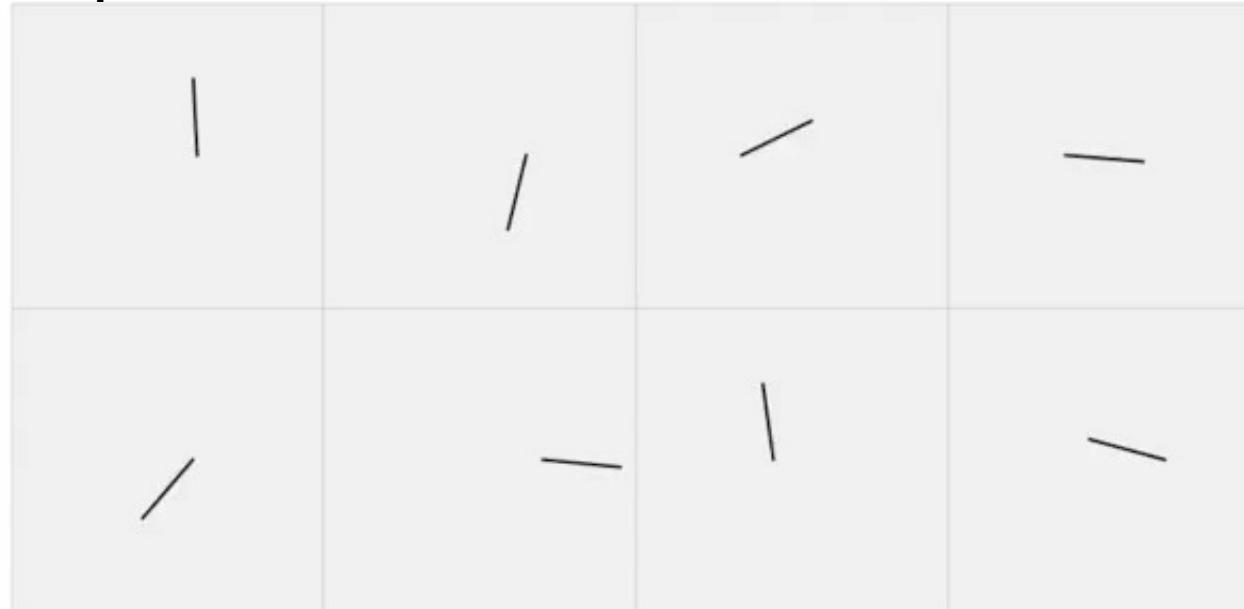
Amortize the problem by learning a latent subspace of optimal solutions

Only search over optimal solutions rather than the entire space

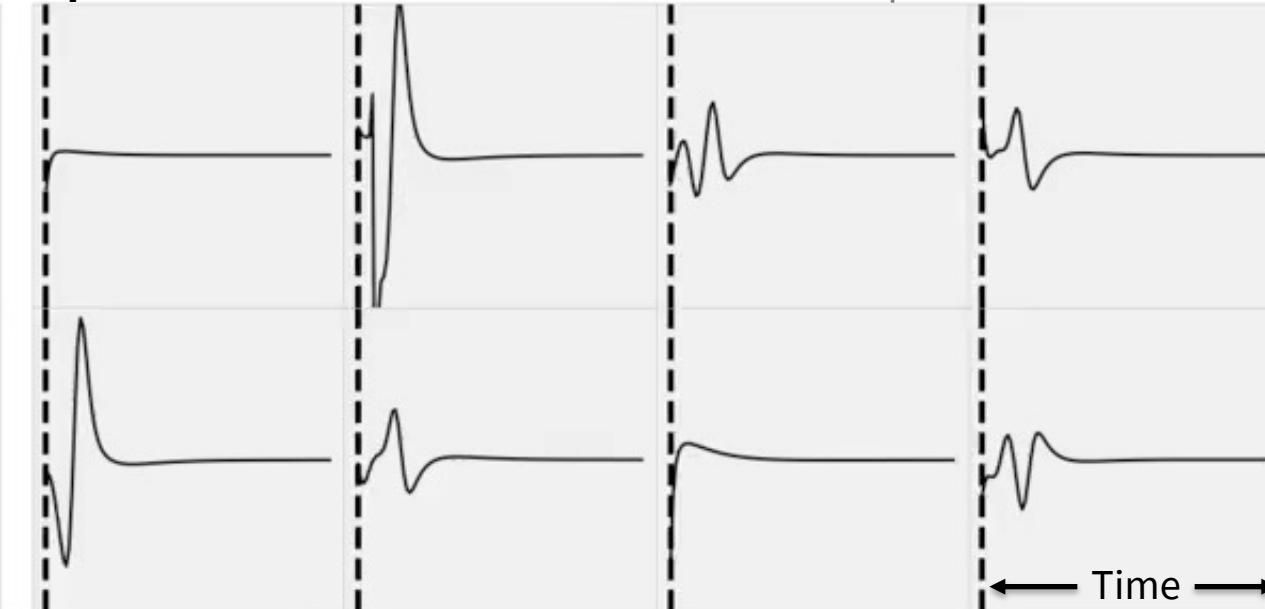
$$x_{1:T}^*, u_{1:T}^* \in \operatorname{argmin}_{x_{1:T}, u_{1:T}} \sum_t \text{cost } C_\theta(x_t, u_t) \text{ s.t. } \begin{array}{l} \text{initial state } x_1 = x_{\text{init}} \\ \text{dynamics } x_{t+1} = f_\theta(x_t, u_t) \\ \text{constraints } u_t \in \mathcal{U} \end{array}$$



Cartpole videos



Optimal controls over time — force on the cartpole



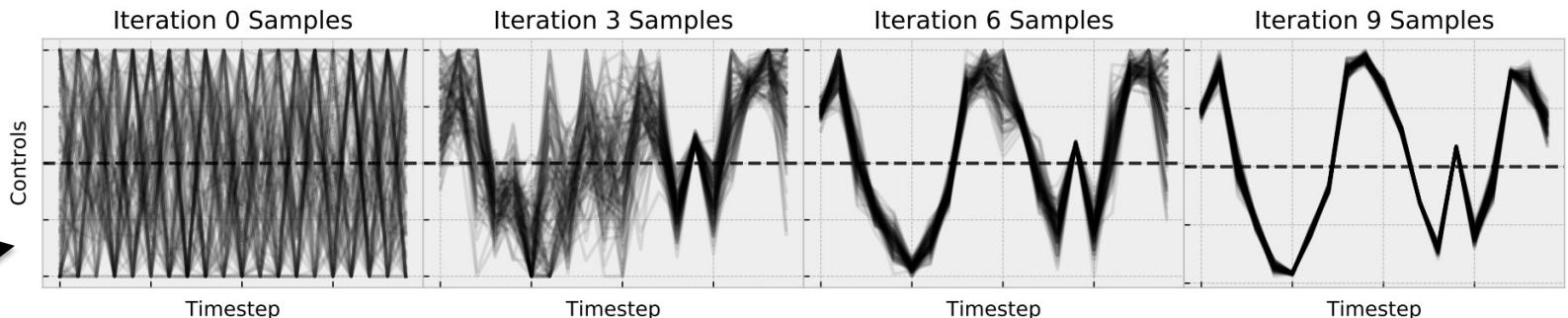
Amortization via learning latent subspaces

The differentiable cross-entropy method. Amos and Yarats, ICML 2020.

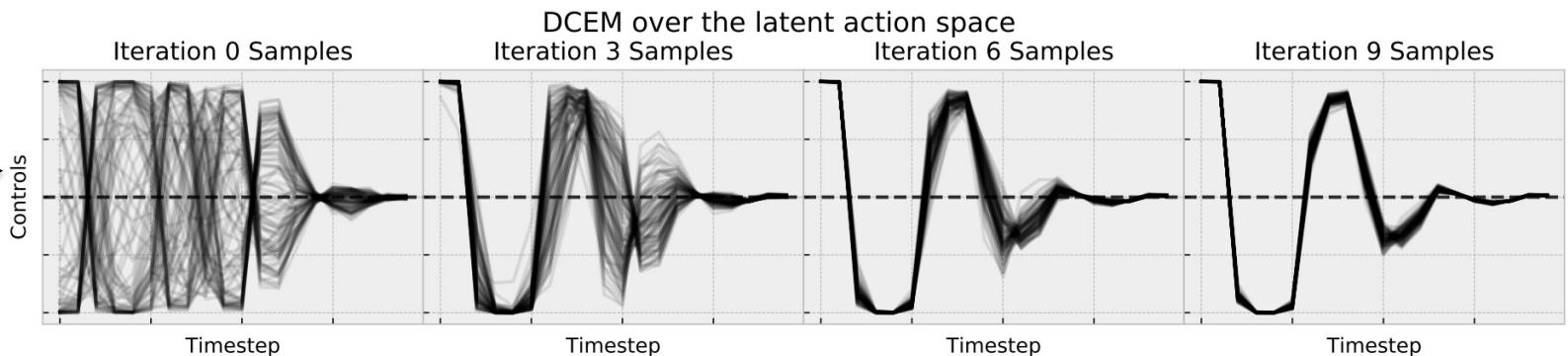
CEM over the full action space

$$u^* = \operatorname{argmin}_{u \in [0,1]^N} f(u)$$

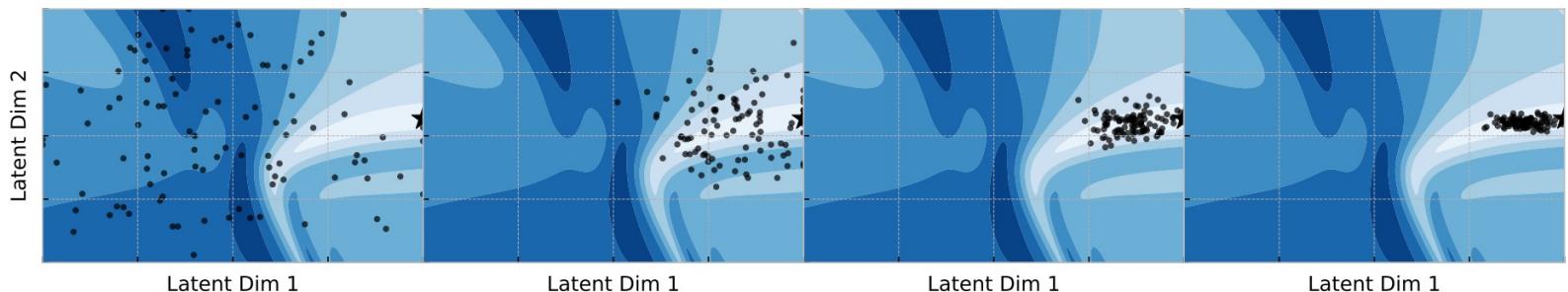
Full control sequence space



Subspace of optimal solutions



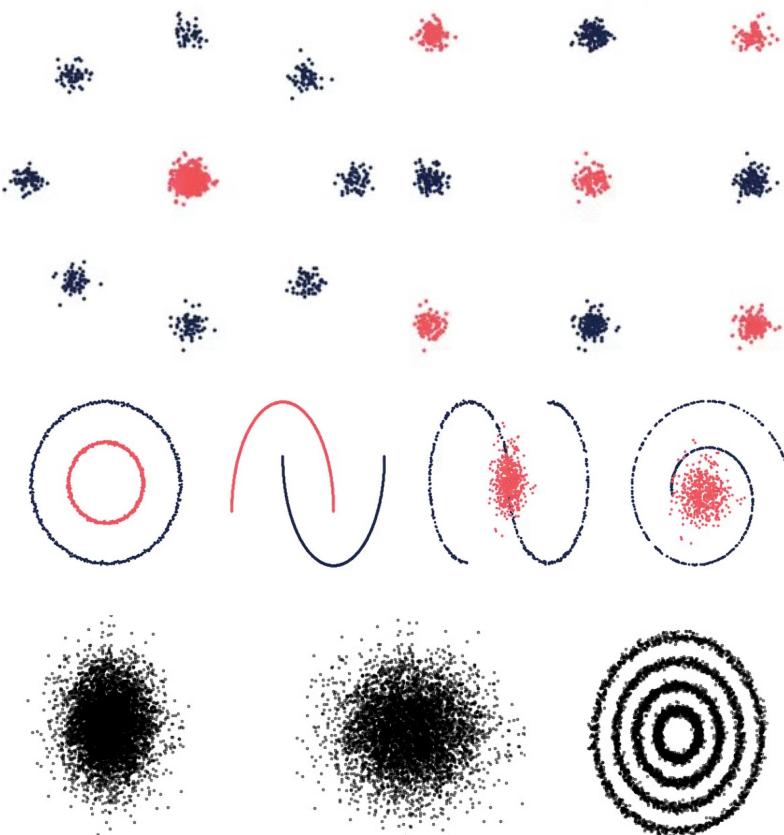
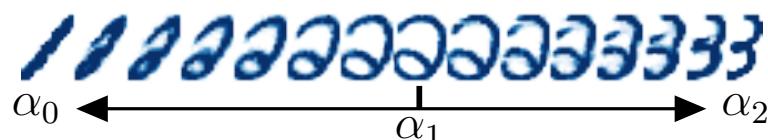
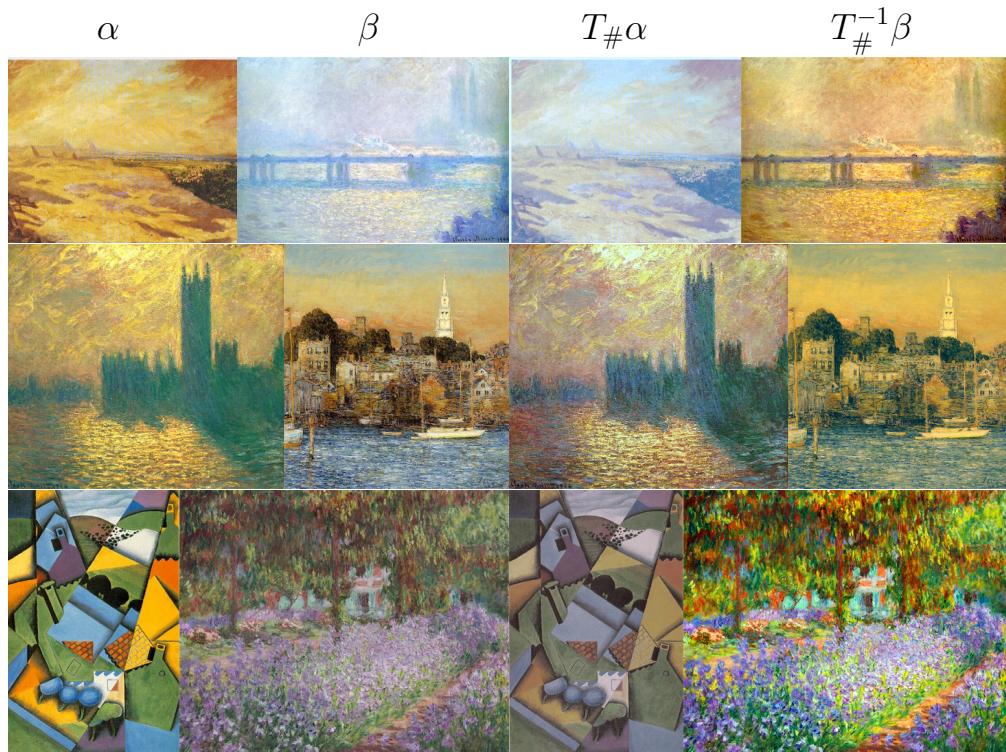
Latent space of optimal solutions



New work: amortization in optimal transport

Meta Optimal Transport. Amos et al., 2022

On amortizing convex conjugates for optimal transport. Amos, ICLR 2023



Closing thoughts

Optimization expresses **non-trivial reasoning operations**

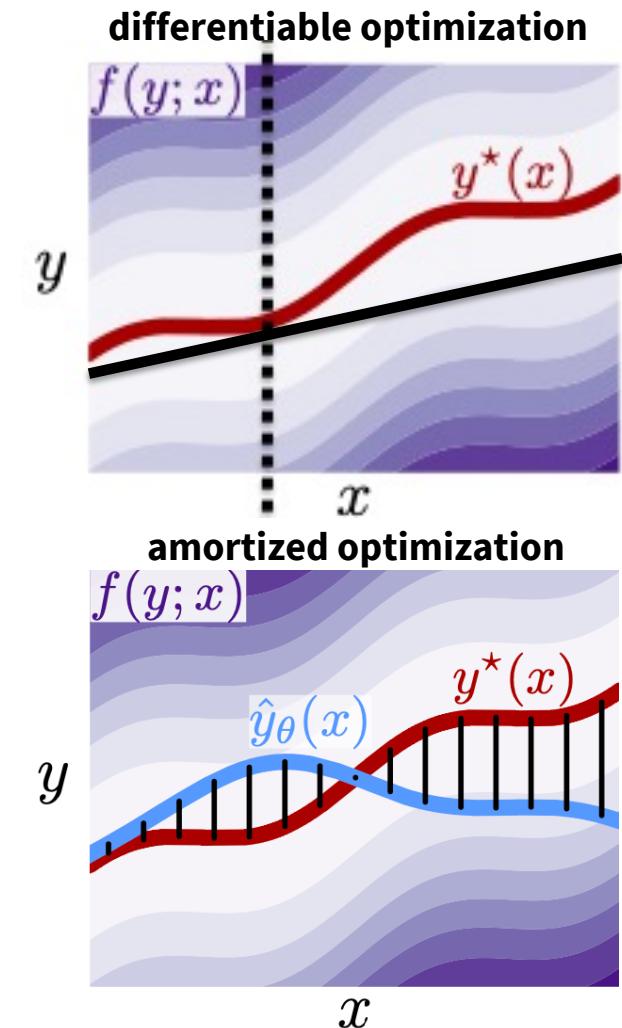
- Integrates nicely with machine learning by **seeing it as a function**

Differentiable optimization is a powerful primitive within larger systems

- Theoretical and engineering foundations are here
- Can be propagated through and learned, just like any layer
- Provides a perspective to analyze existing models and layers

Amortized optimization enables fast learning-augmented solvers

- Foundations are likewise here
- Provides a perspective on existing methods in RL, VI, OT



Future directions

Goal: build intelligent systems that **understand and interact** with the world

Why? To advance scientific and engineering discoveries

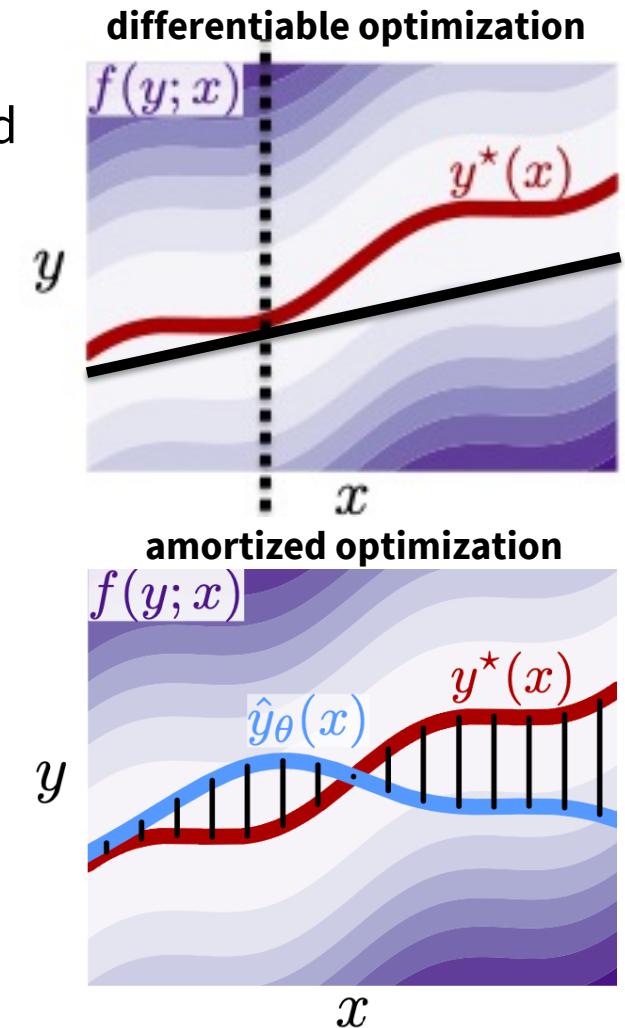
Advancing **optimization** and **machine learning foundations** is crucial

Differentiable optimization — $\frac{\partial}{\partial x} y^*(x)$

- *Differentiable control: learning control-aware representations*
- *Numerics: “differentiation”, ill-behaved functions, approximations*
- *Generalizing to non-convex, non-Euclidean, and/or discrete settings*

Amortized optimization — $\hat{y}_\theta(x) \approx y^*(x)$

- *Numerics: amortizing with distributions and Markov kernels*
- *Building general learning-augmented optimizers*
- *Amortization models that use differentiable optimization*



Learning with differentiable and amortized optimization

Brandon Amos • Meta AI (FAIR) NYC



<http://github.com/bamos/presentations>

[ICML 2017] [Differentiable QPs: OptNet](#)

[NeurIPS 2017] [Differentiable Task-based Model Learning](#)

[NeurIPS 2018] [Differentiable MPC for End-to-end Planning and Control](#)

[NeurIPS 2019] [Differentiable Convex Optimization Layers](#)

[Ph.D. Thesis 2019] [Differentiable Optimization-Based Modeling for ML](#)

[arXiv 2019] [Differentiable Top-k and Multi-Label Projection](#)

[arXiv 2019] [Generalized Inner Loop Meta-Learning](#)

[ICML 2020] [Differentiable Cross-Entropy Method](#)

[ICML 2021] [Differentiable Combinatorial Optimization: CombOptNet](#)

[L4DC 2021] [On the model-based stochastic value gradient](#)

[NeurIPS 2022] [Theseus: Differentiable Nonlinear Optimization](#)

[arXiv 2022] [Meta Optimal Transport](#)

[ICLR 2023] [On amortizing convex conjugates for optimal transport](#)

[Foundations and Trends in ML, to appear] [Tutorial on amortized optimization](#)

Collaborators: Akshay Agrawal, Alexander Rives, Andrew Gordon Wilson, Anselm Paulus, Arnaud Fickinger, Artem Molchanov, Austin Wang, Byron Boots, Caroline Chen, Daniel DeTone, Denis Yarats, Douwe Kiela, Edward Grefenstette, Franziska Meier, Georg Martius, Giulia Luise, Hengyuan Hu, Ievgen Redko, Ivan Jimenez, Jacob Sacks, Jason Liu, Jing Dong, Joseph Ortiz, Joshua Meier, Kyunghyun Cho, Luis Pineda, Maurizio Monge, Michal Rolínek, Mustafa Mukadam, Nathan Lambert, Noam Brown, Omry Yadan, Paloma Sodhi, Phu Mon Htut, Priya Donti, Ricky Chen, Robert Verkuil, Roberto Calandra, Samuel Cohen, Samuel Stanton, Shane Barratt, Shobha Venkataraman, Soumith Chintala, Stephen Boyd, Steven Diamond, Stuart Anderson, Stuart Russell, Taosha Fan, Tom Sercu, Vít Musil, Yann LeCun, Zeming Lin, Zico Kolter