Surrogate Objectives for Batch Policy Optimization in One-step Decision Making

Minmin Chen¹, Ramki Gummadi¹, Chris Harris¹, Dale Schuurmans^{1,2}

COST SENSITIVE CLASSIFICATION

Batch policy optimization

Assume given **complete** data

	a_1	a 2	•••	a _n
<i>x</i> ₁	<i>r</i> ₁₁	<i>r</i> ₁₂	<i>r</i> ₁	<i>r</i> _{1<i>n</i>}
<i>x</i> ₂	<i>r</i> ₂₁	<i>r</i> ₂₂	<i>r</i> ₂	r _{2n}
<i>X</i> 3	<i>r</i> ₃₁	<i>r</i> ₃₂	<i>r</i> ₃	r _{3n}
<i>X</i> 4	<i>r</i> ₄₁	<i>r</i> ₄₂	<i>r</i> 4	r _{4n}
<i>X</i> 5	r_{51}	<i>r</i> ₅₂	<i>r</i> ₅	r _{5n}
<i>x</i> 6	r ₆₁	<i>r</i> ₆₂	<i>r</i> ₆	r _{6n}
	$r_{:1}$	<i>r</i> :2	r _:	r _{:n}
<i>x</i> _m	<i>r_{m1}</i>	<i>r</i> _{m2}	r _m	r _{mn}

Target objective

• expected reward: $\max \sum_{i} \mathbf{r}_{i} \cdot \boldsymbol{\pi}(x_{i})$

Done, right? Not so fast ...

This objective has serious problems

- actually trying to solve: $\max \sum_{i} \mathbf{r}_{i} \cdot \mathbf{f}(q(x_{i}))$
- plateaus everywhere
- can have **exponentially many** local maxima
- nearly impossible to reach a global optima

lso: you already know not to train this way!

to maximize expected reward on **test** contexts

Optimize policy $\pi: X \to \Delta^n$

 $\pi(a \mid x) = e^{q(x)_a - F(q(x))}$ $F(q(x)) = \log \sum_{a} e^{q(x)_{a}}$

 $q: X \to \Re^n$ neural network

Recall: supervised classification

Special case: supervised classification • expected accuracy: $\max \sum_{i} \mathbf{r}_{i} \cdot \boldsymbol{\pi}(x_{i})$

	a_1	<i>a</i> ₂	• • •	a _n
<i>x</i> ₁	0	1	0	0
<i>x</i> ₂	0	0	0	1
<i>X</i> 3	1	0	0	0
<i>X</i> 4	0	0	1	0
<i>X</i> 5	1	0	0	0
<i>x</i> 6	0	1	0	0
1	0	0	1	0
x _m	0	0	0	1

Target objective

But you have never trained with this objective Instead, you used a surrogate objective

maximum likelihood $\max \sum_{i} \mathbf{r}_{i} \cdot \log \boldsymbol{\pi}(x_{i})$

What's going on?

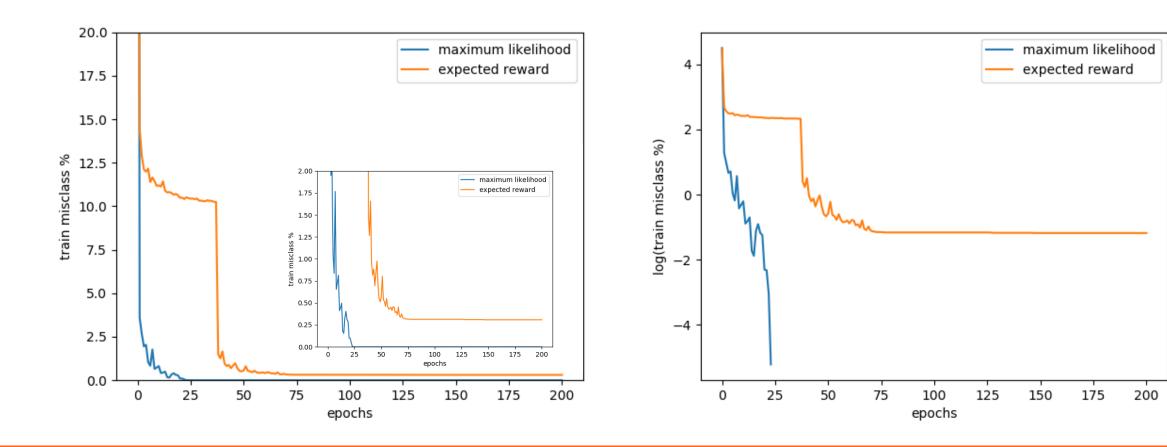
- $\mathbf{r}_i \cdot \boldsymbol{\pi}(x_i)$ is differentiable, that's not the issue
- training with $\mathbf{r}_i \cdot \log \boldsymbol{\pi}(x_i)$ actually achieves better values of $\mathbf{r}_i \cdot \boldsymbol{\pi}(x_i)$ on the training data

Useful properties of maximum likelihood

- $\mathbf{r}_i \cdot \log \boldsymbol{\pi}(x_i)$ is concave in $\mathbf{q}(x_i)$
- it is also calibrated w.r.t. $\mathbf{r}_i \cdot \boldsymbol{\pi}(x_i)$:
- $\forall \epsilon > 0 \exists \delta > 0 \ \mathbf{r} \cdot \log \pi^* \mathbf{r} \cdot \log \pi < \delta \Rightarrow \mathbf{r} \cdot \pi^* \mathbf{r} \cdot \pi < \epsilon$

Target vs surrogate optimization

Misclassification error on MNIST training data



COST SENSITIVE CLASSIFICATION

Definitions

- ▶ Data $\mathcal{D} = \{(x_i, r_i)\}_{i=1}^T$, where $r_i \in \mathbb{R}^K$ specifies reward for each action in context x_i
- ► True risk of a policy is $\mathcal{R}({m \pi}) = -\mathbb{E}[{m \pi}(x) \cdot {m r}]$
- Empirical risk on data set \mathcal{D} is $\hat{\mathcal{R}}(\boldsymbol{\pi},\mathcal{D}) = -\frac{1}{T}\sum_{(x_i,r_i)\in\mathcal{D}} \boldsymbol{\pi}(x_i)\cdot \boldsymbol{r}_i$

Note policies normally represented with composition $\pi(x) = f(q(x))$ where $f(oldsymbol{q})=e^{oldsymbol{q}-oldsymbol{F}(oldsymbol{q})}$ with $F(oldsymbol{q})=\log(oldsymbol{1}\cdot e^{oldsymbol{q}})$

Theorem Even for a linear model $q(x) = W\phi(x)$, the function $r \cdot f(q(x))$ can have exponentially many local maxima in W

Motivation To get around this problem, need to consider *surrogate* training objectives

Calibrated convex surrogate

Definitions

Minimal risk is

$$\mathcal{R}^*(\boldsymbol{r},x) = \inf_{\boldsymbol{\pi}\in\mathcal{P}} \mathcal{R}(\boldsymbol{\pi},\boldsymbol{r},x) = \inf_{\boldsymbol{q}\in\mathcal{Q}} \mathcal{R}(\boldsymbol{f}\circ\boldsymbol{q},\boldsymbol{r},x)$$

► Loss $L^*(r, x) = \inf_{q \in Q} L(q, r, x)$ is *calibrated* w.r.t. \mathcal{R} if: \exists function $\delta(\epsilon, x) \geq 0$ s.t. $\forall \epsilon > 0, x \in X, r \in \mathbb{R}^{K}, q \in Q$:

$$L(\boldsymbol{q},\boldsymbol{r},\boldsymbol{x}) - L^*(\boldsymbol{r},\boldsymbol{x}) < \delta(\epsilon,\boldsymbol{x}) \Rightarrow \mathcal{R}(\boldsymbol{f} \circ \boldsymbol{q},\boldsymbol{r},\boldsymbol{x}) < \mathcal{R}^*(\boldsymbol{r},\boldsymbol{x}) + \epsilon$$

Smoothed risk is

$$S_{\tau}(\boldsymbol{\pi}, \boldsymbol{r}, \boldsymbol{x}) = -\boldsymbol{r} \cdot \boldsymbol{\pi}(\boldsymbol{x}) + \tau \boldsymbol{\pi}(\boldsymbol{x}) \cdot \log \boldsymbol{\pi}(\boldsymbol{x})$$

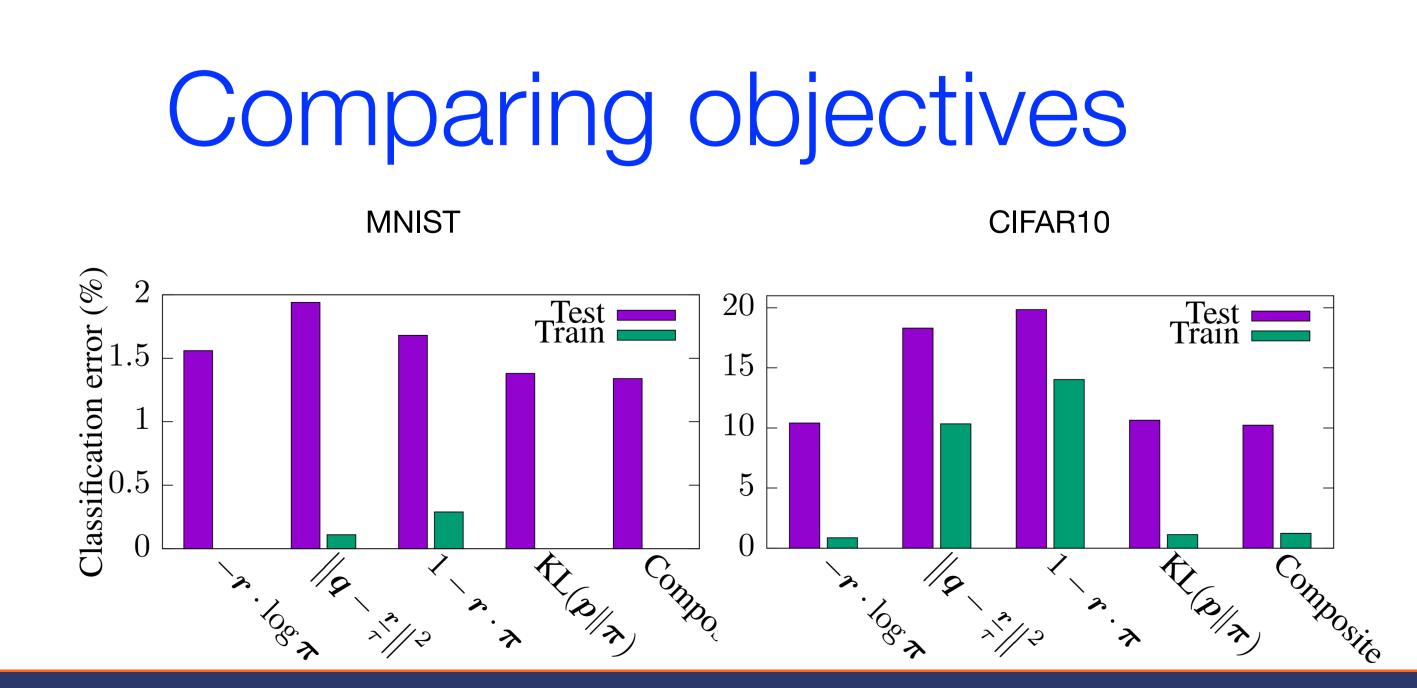
 $\blacktriangleright \ \mathsf{Let} \ \tilde{\boldsymbol{\pi}}_{\tau} = \mathsf{arg} \min_{\boldsymbol{\pi} \in \mathcal{P}} \mathcal{S}_{\tau}(\boldsymbol{\pi})$

Proposition $au < \epsilon / \log K$ implies $\mathcal{R}(ilde{\pi}_{ au}) < \mathcal{R}^* + \epsilon$

Proposition Local smoothed risk is equivalent to $\mathcal{S}_{\tau}(\boldsymbol{\pi}, \boldsymbol{r}, \boldsymbol{x}) = -\tau F(\frac{r}{\tau}) + \tau D_F(\frac{r}{\tau} \| \boldsymbol{q}(\boldsymbol{x}))$

Theorem The surrogate objective $L(\boldsymbol{q}, \boldsymbol{r}, \boldsymbol{x}) = au D_F \left(\boldsymbol{q}(\boldsymbol{x}) + rac{v}{ au} \| rac{r}{ au}
ight) + rac{\tau}{4} \left\| \boldsymbol{q}(\boldsymbol{x}) - rac{r-v}{ au} \right\|^2$ is strongly convex in \boldsymbol{q} and calibrated w.r.t. $\mathcal{S}_{\tau}(\boldsymbol{f} \circ \boldsymbol{q}, \boldsymbol{r} - \boldsymbol{v}, \boldsymbol{x})$ with $\delta(\epsilon, \boldsymbol{x}) = \epsilon$

Experimental evaluation



BATCH CONTEXTUAL BANDITS

Coping with missing data

	a_1	a 2	•••	a _n
<i>x</i> ₁	0	$\frac{r_1}{\beta_1}$	0	0
<i>x</i> ₂	0	0	0	$\frac{r_2}{\beta_2}$
x ₂ x ₃ x4	$\frac{r_3}{\beta_3}$	0	0	0
<i>X</i> 4	0	0	$\frac{r_4}{\beta_4}$	0
×5 ×6	$rac{r_5}{eta_5}$	0	0	0
<i>x</i> 6	0	$\frac{r_6}{\beta_6}$	0	0
1	0	0	$\frac{r_{\cdot}}{\beta_{\cdot}}$	0
x _m	0	0	0	$\frac{r_m}{\beta_m}$

Optimize policy $\pi: X \to \Delta^n$

Example portance corrected expected reward $\max \sum \frac{\pi(a_i | x_i)}{\beta_i} r_i$

where β are proposal probabilities from behavior strateg

We already know this is a poor objective but what about missing data inference?

Equivalent to $\max \hat{\mathbf{r}} \cdot \boldsymbol{\pi}$ using $\hat{\mathbf{r}}_i = \mathbf{1}_{a_i \frac{r_i}{\beta}}$

That is

exaggerate observed values by 1/

• fill in all unobserved values with 0

This is a pretty lame inference principle But ... its unbiased! $\mathbb{E}[\hat{\mathbf{r}} | x] = \sum_{a} \beta_{a} \mathbf{1}_{a} \frac{r_{a}}{\beta} = \sum_{a} \mathbf{1}_{a} r_{a} = \mathbf{r}$

Missing data inference

	<i>a</i> 1	<i>a</i> 2	•••	a _n
<i>x</i> ₁	$ au q_{11}$	$\tau q_{12} + \lambda (r_1 - \tau q_{12})$	$ au q_{1}$	$ au q_{1n}$
<i>x</i> ₂	$ au q_{21}$	τq_{22}	$ au q_{2}$	$\tau q_{2n} + \lambda (r_2 - \tau q_{2n})$
<i>x</i> 3	$\tau q_{31} + \lambda (r_3 - \tau q_{31})$	τ q 32	<i>τq</i> ₃	$ au q_{3n}$
<i>x</i> 4	$ au q_{41}$	$ au q_{42}$	$\tau q_{4} + \lambda (r_4 - \tau q_{4})$	$ au q_{4n}$
<i>X</i> 5	$ au q_{51} + \lambda (r_5 - au q_{51})$	τ q 52	$ au q_{5}$	$ au q_{5n}$
<i>x</i> 6	$ au q_{61}$	$\tau q_{62} + \lambda (r_6 - \tau q_{62})$	$ au q_{6}$	$ au q_{6n}$
	$ au oldsymbol{q}_{:1}$	$ au q_{:2}$	$ au \mathbf{q}_{\ldots} + \lambda (\mathbf{r}_{\cdots} au \mathbf{q}_{\ldots})$	$ au q_{:n}$
x _m	$ au q_{m1}$	$ au q_{m2}$	$ au q_{m}$	$ au q_{mn} + \lambda (r_m - au q_{mn})$

Improvement

"douby robust estimation" instead of filling in with 0s

• fill in with guesses from a model $\mathbf{q}(x)$

 $\hat{\mathbf{r}} = \tau \mathbf{q} + \lambda \mathbf{1}_a (r - \tau q_a)$ Also unbiased

• as long as $\lambda = 1/\beta_i$ out still alters observed data

Where should the model come from?

- could use a separate critic
- train via least squares, then optimize π • works okay, but not great

there is only one action value function for single-step decision making, r(x, a)actor-critic approaches trivialized

Unified approach

Unified approach

- actor and critic are same model
- $\boldsymbol{\pi} = e^{\mathbf{q} F(\mathbf{q})}$ where $F(\mathbf{q}) = \log \mathbf{1} \cdot e^{\mathbf{q}}$
- use logits $\tau \mathbf{q}(x)$ to predict rewards

$$q(x,a) \approx \frac{r(x,a)}{\tau}$$

Can combine with previous objectives • KL($\pi \| \hat{\mathbf{p}}$) where $\hat{\mathbf{p}} = e^{\hat{\mathbf{r}}/\tau - F(\hat{\mathbf{r}}/\tau)}$

• KI $(\hat{\mathbf{n}} \| \boldsymbol{\pi})$

$$\mathsf{KL}(\boldsymbol{\pi} \| \hat{\mathbf{p}}) \leq \mathsf{KL}(\hat{\mathbf{p}} \| \boldsymbol{\pi}) + \frac{\tau}{4} \| \hat{\mathbf{r}} / \tau - \mathbf{q} \|^2$$

these are somewhat sensitive to ranking unlike least squares

Empirical Bayes estimation

- optimize hyperparameters q
- (neural network)
- integrate out parameters ξ

Example marginal likelihood

 $-\log p(r_0 \,|\, a_0, \mathbf{q})$ $= -\log \int p(r_0 | a_0, \xi) p(\xi | \mathbf{q}) d\xi$ $= \frac{1}{2\sigma^2} (\phi(a_0) \cdot q - r_0)^2 + \frac{1}{2} \log \sigma^2 + c$

essentially least squares regression

Can alternatively use surrogates nin **KL**(prior || posterior) min **KL**(posterior || prior) $\approx \min I(\xi; r_0)$



BATCH CONTEXTUAL BANDITS

Reward estimation For x, a, r_a , parameters $\lambda(x, a)$, τ , estimate full $\hat{\boldsymbol{r}}(\boldsymbol{x}) = au \boldsymbol{q}(\boldsymbol{x}) + \mathbf{1}_a \lambda(\boldsymbol{x}, \boldsymbol{a})(\boldsymbol{r}_a - au \boldsymbol{q}(\boldsymbol{x})_a)$

Surrogate objective

Definition Optimal imputed local risk and suboptimality gap $\mathcal{S}^*_{ au}(\hat{\pmb{r}}, \pmb{x}) = \inf_{\pmb{q} \in \mathcal{Q}} \mathcal{S}_{ au}(\pmb{f} \circ \pmb{q}, \hat{\pmb{r}}, \pmb{x}), \; \mathcal{G}_{ au}(\pmb{\pi}, \hat{\pmb{r}}, \pmb{x}) = \mathcal{S}_{ au}(\pmb{\pi}, \hat{\pmb{r}}, \pmb{x}) - \mathcal{S}^*_{ au}(\hat{\pmb{r}}, \pmb{x})$

Proposition $\forall \boldsymbol{q}, \tau > 0, (x, a, r_a) : \tau D_F(\frac{\hat{r}(x)}{\tau} \| \boldsymbol{q}(x)) = \mathcal{G}_{\tau}(\boldsymbol{f} \circ \boldsymbol{q}, \hat{\boldsymbol{r}}, x)$

$$\begin{array}{ll} \textbf{Theorem} & \forall \boldsymbol{q}, \ \tau > 0, \ (x, a, r_a), \ v: \\ & L(\boldsymbol{q}, \hat{\boldsymbol{r}}, x) \geq \tau D_F \left(\frac{\hat{\boldsymbol{r}}(x)}{\tau} \middle\| \boldsymbol{q}(x) + \frac{\boldsymbol{v}}{\tau} \right) = \mathcal{G}_{\tau}(\boldsymbol{f} \circ \boldsymbol{q}, \hat{\boldsymbol{r}}, x) \geq 0 \\ & L \ \text{calibrated w.r.t.} \ \mathcal{S}_{\tau}(\boldsymbol{f} \circ \boldsymbol{q}, \hat{\boldsymbol{r}} - \boldsymbol{v}, x) \end{array}$$

Optimization Given $\mathcal{D} = \{(x_i, a_i, r_i, \beta_i)\}$, context, act, reward, prob $\min_{\boldsymbol{q}\in\mathcal{Q}} \hat{L}(\boldsymbol{q},\mathcal{D}) \quad \text{where} \quad \hat{L}(\boldsymbol{q},\mathcal{D}) = \frac{1}{T} \sum_{(x_i,a_i,r_i,\beta_i)\in\mathcal{D}} L(\boldsymbol{q},\hat{\boldsymbol{r}},x_i)$

Analysis

Definition Expected smoothed risk quantities we seek to control: $\mathcal{S}_{ au}(oldsymbol{\pi}) = \mathbb{E}[\mathcal{S}_{ au}(oldsymbol{\pi},oldsymbol{r},oldsymbol{x})], \mathcal{S}^*_{ au} = \mathsf{inf}_{oldsymbol{q}\in\mathcal{Q}}\,\mathcal{S}_{ au}(oldsymbol{f}\circoldsymbol{q}), \mathcal{G}_{ au}(oldsymbol{\pi}) = \mathcal{S}_{ au}(oldsymbol{\pi}) - \mathcal{S}^*_{ au}$

- **Theorem** For any q, \hat{r} such that $\mathbb{E}[\hat{r}|x] = \mathbb{E}[r|x]$, and baseline v: $\mathbb{E}[L(\boldsymbol{q}, \hat{\boldsymbol{r}}, \boldsymbol{x})] \geq \mathbb{E}\left[\tau D_F\left(\frac{\hat{\boldsymbol{r}}(\boldsymbol{x})}{\tau} \middle\| \boldsymbol{q}(\boldsymbol{x}) + \frac{\boldsymbol{v}}{\tau}\right)\right] \geq \mathcal{G}_{\tau}(\boldsymbol{f} \circ \boldsymbol{q}) \geq 0.$
- **Lemma** $\forall \tau, \delta > 0 \exists$ constant *C* s.t. w.p. at least 1δ : $\mathbb{E}\left[D_F\left(rac{\hat{r}(x)}{ au} \| oldsymbol{q}(x)
 ight)
 ight] \leq \hat{D}_F(oldsymbol{q},\mathcal{D}) + rac{\mathcal{C}}{\sqrt{ au}} \quad orall oldsymbol{q} \in \mathcal{H}.$
- **Theorem** $\forall v, \tau, \delta > 0, \exists C \text{ s.t. w.p. at least } 1 \delta$: if $\hat{L}(\boldsymbol{q},\mathcal{D}) < rac{ au C}{\sqrt{T}}$ for $\boldsymbol{q} \in \mathcal{H}$ then $\mathcal{G}_{ au}(\boldsymbol{f} \circ \boldsymbol{q}) \leq rac{2 au C}{\sqrt{T}}$

Experimental evaluation

