

Examples are not Enough, Learn to Criticize! Criticism for Interpretability

Been Kim, Rajiv Khanna, Oluwasanmi Koyejo

Wittawat Jitkrittum

Gatsby Machine Learning Journal Club

16 Jan 2017

Summary

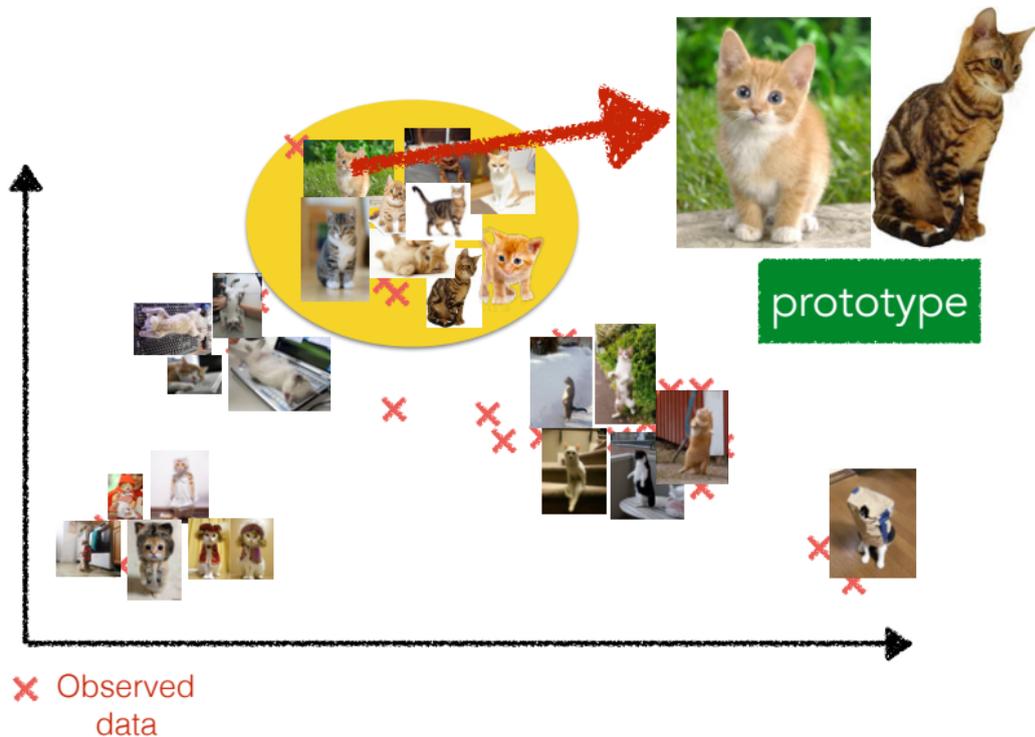
Examples are not Enough, Learn to Criticize! Criticism for Interpretability
Been Kim, Rajiv Khanna, Oluwasanmi O. Koyejo
NIPS 2016.

Given a big dataset, want to do 2 things:

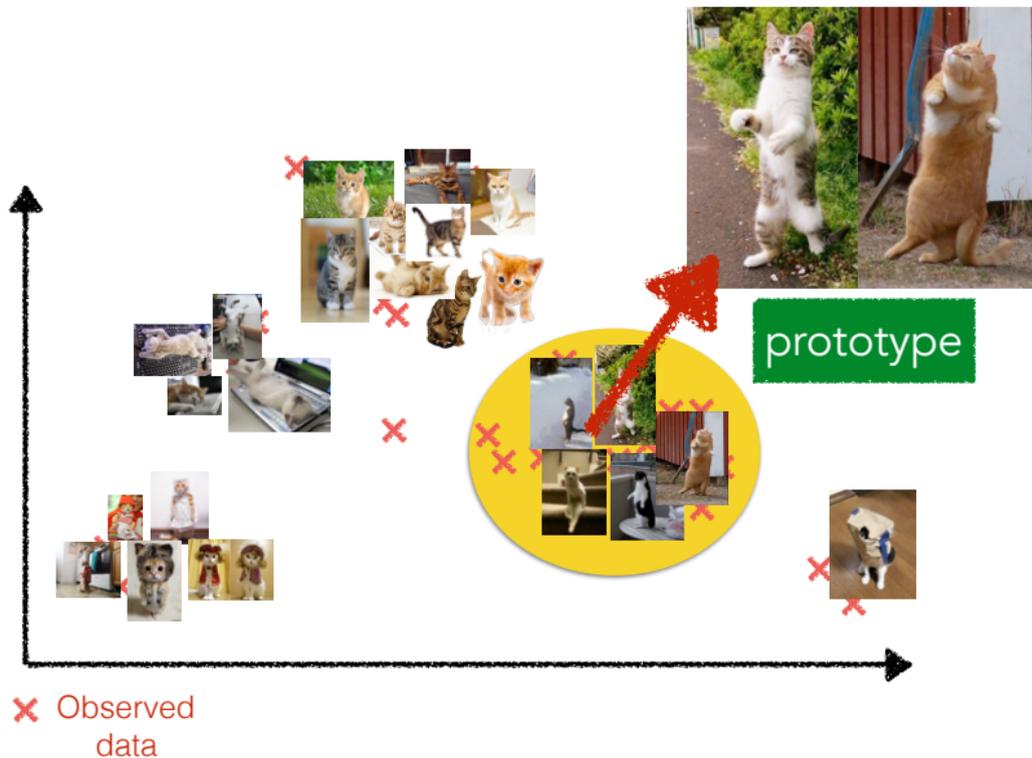
- 1 **Summarize:** Find typical examples = prototypes. Majorities.
 - 2 **Criticize:** Find atypical examples that are not covered by the prototypes. Minorities.
- Many existing works focus on only [1] e.g., K-medoid, set cover.
 - Main message: [2] is also important.
 - Use kernel MMD as the objective.

(Some slides are stolen from Been Kim.)

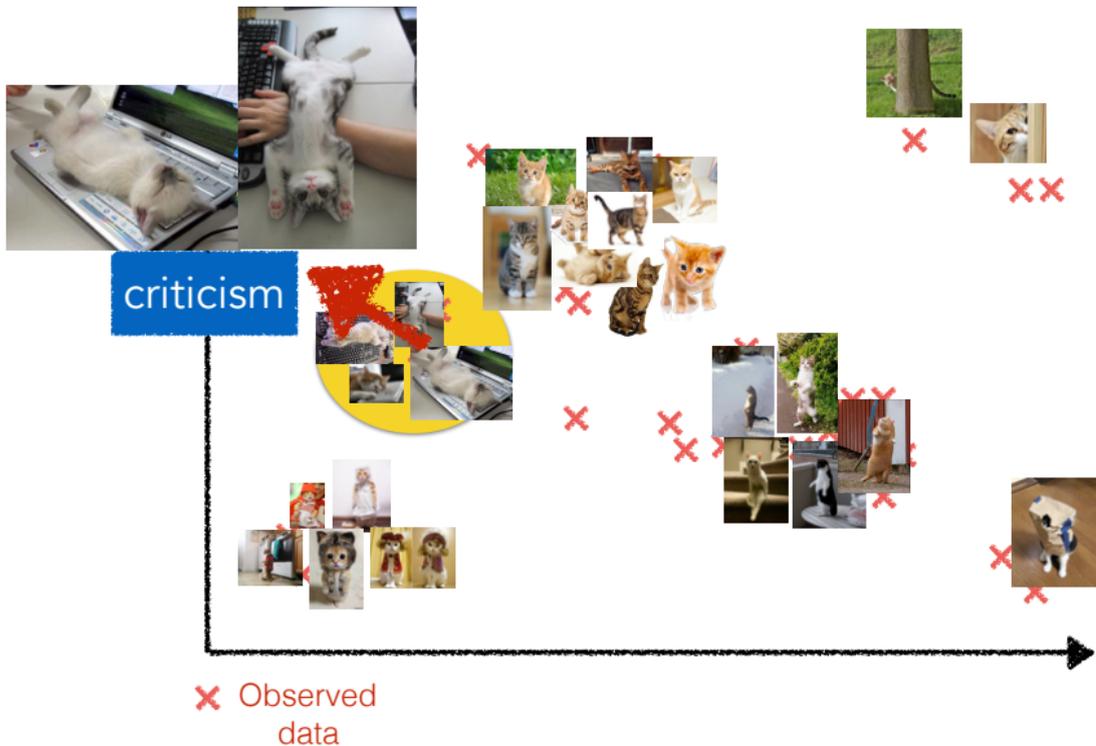
Understanding data through examples



Understanding data through examples



Understanding data through examples



Maximum Mean Discrepancy (MMD)

- k : a kernel associated with RKHS \mathcal{H} s.t. $k(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}}$.
- Two sets of samples: $X = \{x_i\}_{i=1}^n \sim P$, $Z = \{z_i\}_{i=1}^m \sim Q$.
- Empirical MMD:

$$\begin{aligned} \text{MMD}^2(X, Z) &= \left\| \frac{1}{n} \sum_{i=1}^n \phi(x_i) - \frac{1}{m} \sum_{j=1}^m \phi(z_j) \right\|_{\mathcal{H}}^2 \\ &= \frac{1}{n^2} \sum_{i,j=1}^n k(x_i, x_j) - \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m k(x_i, z_j) + \frac{1}{m^2} \sum_{i,j=1}^m k(z_i, z_j). \end{aligned}$$

- **Summarization:** Choose subset indices $S \subset \{1, \dots, n\}$ to minimize $\text{MMD}^2(X, X_S)$. (Cf. kernel herding).
 - Pick $|S| = m$ points to preserve the moments as defined by $\phi(\cdot)$.

Proposal: MMD-critic for Prototypes

- Given $X = \{x_i\}_{i=1}^n$, define a maximization objective $J_b(S)$

$$\begin{aligned} J_b(S) &= \frac{1}{n^2} \sum_{i,j=1}^n k(x_i, x_j) - \text{MMD}^2(X, X_S) \\ &= \underbrace{\frac{2}{n|S|} \sum_{i=1}^n \sum_{j \in S} k(x_i, x_j)}_{\text{relevancy}} - \underbrace{\frac{1}{|S|^2} \sum_{i \in S} \sum_{j \in S} k(x_i, x_j)}_{\text{redundancy}} \end{aligned}$$

- $\frac{1}{n^2} \sum_{i,j=1}^n k(x_i, x_j)$ is constant. Added so that $J_b(\emptyset) = 0$ (“normalized”).
- Select m_* prototypes by (discrete optimization)

$$\max_{S \subset \{1, \dots, n\}, |S| \leq m_*} J_b(S).$$

Optimization Guarantees

- Def: $F(S)$ is **normalized** if $F(\emptyset) = 0$.
- Def: $F(S)$ is **monotonic** if $U \subseteq V \subseteq \{1, \dots, n\}$ implies $F(U) \leq F(V)$.
- Def: $F(S)$ is **submodular** if for all $U, V \subseteq \{1, \dots, n\}$,

$$F(U \cup V) + F(U \cap V) \leq F(U) + F(V).$$

- Will show that $J_b(S)$ is monotonic, submodular under some conditions.
- Then, use greedy forward search. At each iteration t ,

$$S_{t+1} = S_t \cup \left\{ \arg \max_{u \in \{1, \dots, n\} \setminus S_t} J_b(S_t \cup \{u\}) \right\}.$$

Theorem (Nemhauser et al. (1978))

If F is normalized, monotonic, submodular, then the greedy approach achieves at least $(1 - e^{-1}) \max_{|S| \leq m_} F(S)$.*

Variational View of MMD

$$\begin{aligned}\text{MMD}(P, Q) &= \left\| \mathbb{E}_{X \sim P}[\phi(X)] - \mathbb{E}_{Y \sim Q}[\phi(Y)] \right\|_{\mathcal{H}} \\ &= \sup_{f \in \mathcal{H}, \|f\| \leq 1} \mathbb{E}_{X \sim P}[f(X)] - \mathbb{E}_{Y \sim Q}[f(Y)].\end{aligned}$$

- $\arg \sup$ is the **witness function**:

$$f(\mathbf{x}) = \mathbb{E}_{X' \sim P}[k(\mathbf{x}, X')] - \mathbb{E}_{Y' \sim Q}[k(\mathbf{x}, Y')].$$

- $f(\mathbf{x}) > 0$ in high density areas of P .
- $f(\mathbf{x}) < 0$ in high density areas of Q .
- Magnitude $|f(\mathbf{x})|$ indicates the density difference at \mathbf{x} .

-
- For our purpose, the empirical witness associated with $\text{MMD}(X, X_S)$:

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n k(\mathbf{x}, x_i) - \frac{1}{|S|} \sum_{j \in S} k(\mathbf{x}, x_j).$$

MMD-critic for Criticisms

- Criticisms of S are points with high magnitude of the witness f

$$C = \arg \max_{C \subseteq \{1, \dots, n\} \setminus S, |C| < c_*} L(C) + \log \det K_{C,C}$$

$$L(C) = \sum_{l \in C} |f(x_l)| = \sum_{l \in C} \left| \frac{1}{n} \sum_{i=1}^n k(x_i, x_l) - \frac{1}{|S|} \sum_{j \in S} k(x_j, x_l) \right|.$$

- Regularizer $\log \det K_{C,C}$ is high when $\{x_l\}_{l \in C}$ are diverse.
- $L(C) + \log \det K_{C,C}$ is sub-modular. Greedy optimization.

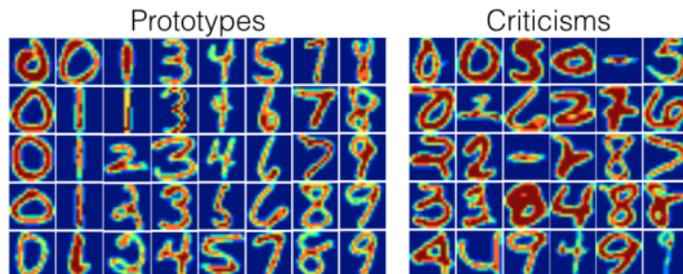
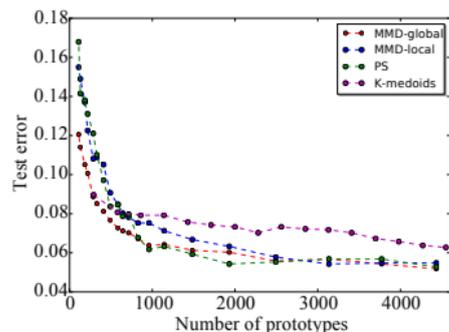
-
- The whole procedure gives summary points S , and criticisms C .

Quality of the Prototypes

- Find prototypes of USPS handwritten digits.
- Gaussian kernel: $k(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$.
- Use 1-NN (nearest prototype) classification error as the quality measure.
- Let $y_i \in \{1, \dots, 10\}$ be the class label of x_i .
- Given \hat{x} , the nearest prototype classifier predicts y_{i^*} , where

$$i^* = \arg \min_{i \in S} \|\phi(\hat{x}) - \phi(x_i)\|_{\mathcal{H}}^2 = \arg \min_{i \in S} k(\hat{x}, x_i).$$

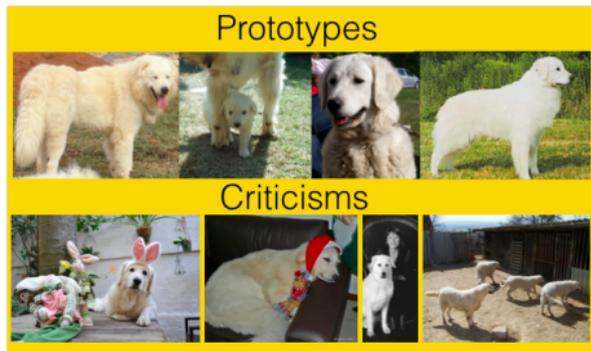
Performance on USPS Data



- MMD-local: Use $\exp(-\gamma\|x_i - x_j\|^2)[y_i = y_j]$. Supervised kernel to find the prototypes.
- MMD-global: Use the usual Gaussian kernel.
- PS: Prototype Selection of Bien and Tibshirani, 2011.
- Features = raw pixels.

Qualitative Measure: Prototype and Criticisms

- Two types of dog breeds from Imagenet.
- Features = image embeddings from He et al., 2015.



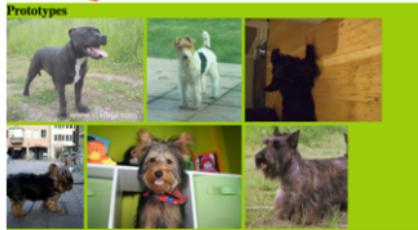
Eval3

Pilot study with human subjects

- Definition of interpretability: A method is interpretable if a user can correctly and efficiently predict the method's results.
- Task: Assign a new data point to one of the groups using 1) all images 2) prototypes 3) prototypes and criticisms 4) small set of randomly selected images



a new data point



group 1



group 2

Eval3

Pilot study with human subjects

- Definition of interpretability: A method is interpretable if a user can correctly and efficiently predict the method's results.
- Task: Assign a new data point to one of the groups using 1) all images 2) prototypes 3) prototypes and criticisms 4) small set of randomly selected images



a new data point



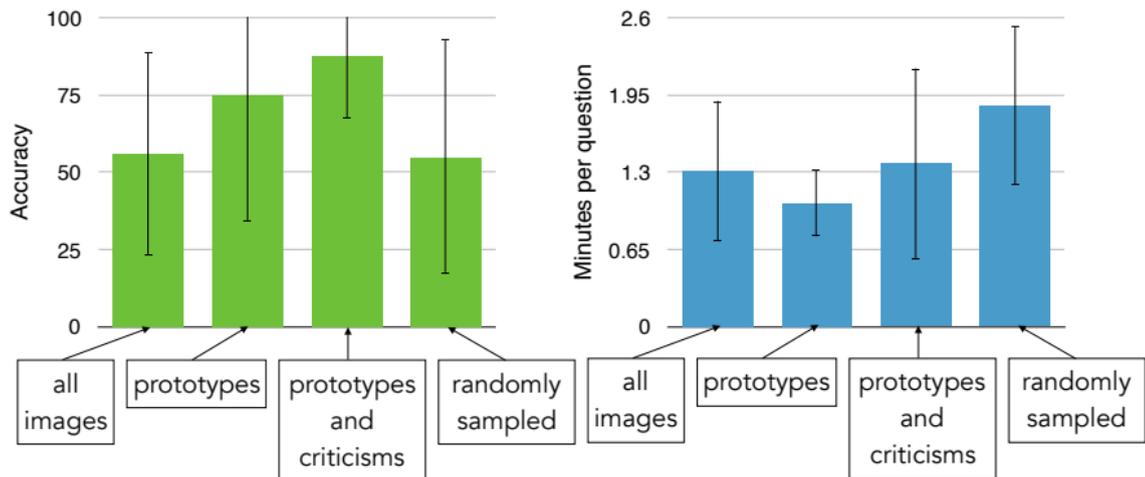
group 1



group 2

Eval3

Pilot study with human subjects



Comment:

"[Proto and Criticism Condition resulted in] less confusion from trying to discover hidden patterns in a ton of images, more clues indicating what features are important"

n = 3

21 questions each

Some Questions

- What happens when there is no $\log \det K_{C,C}$?
- Quantify the effect of the image embeddings from He et al., 2015.
- There are only 3-4 human subjects.
- Possible to do a continuous optimization without selecting a subset?

Lemma 1: $J_b(S)$ Is Linear in K

- Let $K \in \mathbb{R}^{n \times n}$ such that $k_{ij} = k(x_i, x_j)$.
- Prototype objective:

$$\begin{aligned} J_b(S) &= \frac{2}{n|S|} \sum_{i=1}^n \sum_{j \in S} k(x_i, x_j) - \frac{1}{|S|^2} \sum_{i \in S} \sum_{j \in S} k(x_i, x_j) \\ &= \frac{2}{n|S|} \sum_{i=1}^n \sum_{j=1}^n [j \in S] k(x_i, x_j) - \frac{1}{|S|^2} \sum_{i=1}^n \sum_{j=1}^n [i \in S][j \in S] k(x_i, x_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n \left(\frac{2}{n|S|} [j \in S] - \frac{1}{|S|^2} [i \in S][j \in S] \right) k_{ij} \\ &= \sum_{i=1}^n \sum_{j=1}^n a_{ij}(S) k_{ij} := \langle A(S), K \rangle, \end{aligned}$$

- Matrix inner product: $\langle A, B \rangle = \sum_i \sum_j a_{ij} b_{ij}$.

Theorem 2.1: Monotone Linear Forms

- Given $H \in \mathbb{R}^{n \times n}$ s.t. $0 \leq h_{ij} \leq h_*$ where $h_* := \max_{i,j} h_{ij} > 0$.
- Define $E \in \{0, 1\}^{n \times n}$ s.t. $e_{ij} = [h_{ij} = h_*]$.
- Define $F(B, S) := \langle A(S), B \rangle$.
- Let $m := |S|$. Define

$$\alpha(n, m) = \frac{F(E, S \cup \{u\}) - F(E, S)}{F(1 - E, S)},$$

for all $u \in S$.

- If for all i, j s.t. $[h_{ij} \neq h_*]$, for all $m \in \{0, \dots, n\}$, $h_{i,j} \leq h_* \alpha(n, m)$, then $F(H, S)$ is monotone.

-
- A similar statement to guarantee that $F(H, S)$ is submodular.

(Corollary) Monotone Submodularity for MMD

Assume

- 1 K is s.t., $k_{ij} \geq 0$.
- 2 $k_{i,i} = k_* > 0$ for all $i \in \{1, \dots, n\}$.
- 3 K is diagonally dominant i.e., $\sum_{j \neq i} |k_{i,j}| < |k_{i,i}|$ for all i .
- 4 $k_{i,j} \leq \frac{k_*}{n^3 + 2n^2 - 2n - 3}$

Then, $J_b(S)$ is monotone submodular.

- For a fixed n , and $k_{i,j} = \exp(-\gamma \|x_i - x_j\|^2)$, there exists γ such that (3), (4) are satisfied.
- What if n is very large?

Questions?

Thank you