

Big Data



① volume

- sample size is big
- feature dimensionality is big

② variety

- multiple formats: social, video, unstructured
- problem is big and has many related tasks

③ velocity

- batch → real-time

Big Data



① volume

- sample size is big
- feature dimensionality is big

② variety

- multiple formats: social, video, unstructured
- problem is big and has many related tasks

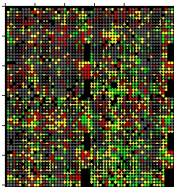
③ velocity

- batch → real-time

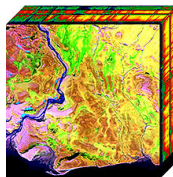
Real-world data are often high-dimensional



text



bioinformatics



hyperspectral



image

- feature extraction / feature selection

- 1 **filter** approach
 - as **preprocessing** step (no interaction with learning algorithm)
- 2 **wrapper** approach
 - use the learning algorithm to score feature subsets
- 3 **embedded** approach
 - perform feature selection and learning **simultaneously** → sparse solution

Regularized risk minimization

minimize $\text{loss } \ell(w) + \text{sparsity-inducing regularizer } \Omega(w)$

Lasso (Tibshirani, 1996)

- $\Omega(w) = \|w\|_1$



- highly correlated features \rightarrow tends to arbitrarily select only one of them

1

- _____



- 10/10/2017

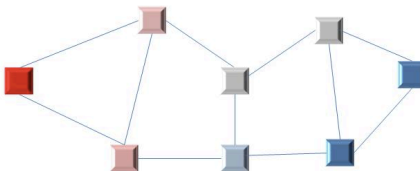


- $$\mathcal{C}(x) = \{x, x+1, \dots, x+d-1\}$$

- $$\Omega(W) = \|W\|_1 + \lambda \sum_{i=1}^{d-1} W_i^{p-1} \quad [W_i = W_{i+1}]$$

Graph lasso

- features are ordered in a graph $G = (V, E)$



- $\Omega(w) = \|w\|_1 + \sum_{(i,j) \in E} |w_i - w_j|$
- encourages coefficients for nearby graph nodes to be similar

Requires the **structure** to be **known** in advance

- group structure / sequential ordering / graph
- may not be available

Sparse Modeling with Automatic Feature Grouping

Goal

Feature coefficients

- 1 sparse
- 2 grouped **automatically**
- 3 have similar magnitudes for features in the same group

Example

- group of dummy variables for the same categorical variable
- protein-protein interaction networks: groups of co-regulated genes
- text classification: groups of correlated words

Advantages

- variance reduction \rightarrow better accuracy
- simpler model \rightarrow better interpretation

Octagonal Shrinkage and Clustering Algorithm for Regression

- ℓ_1 -regularizer: encourages sparsity
- (convex) pairwise ℓ_∞ -regularizer: tries to tie every coefficient pair $|w_i|, |w_j|$ together



0001

- “ ”

- 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2632, 2633, 2634, 2635, 2636, 2637, 2638, 2639, 2640, 2641, 2642, 2643, 2644, 2645, 2646, 2647, 2648, 2649, 2650, 2651, 2652, 2653, 2654, 2655, 2656, 2657, 2658, 2659, 2660, 2661, 2662, 2663, 2664, 2665, 2666, 2667, 2668, 2669, 2670, 2671, 2672, 2673, 2674, 2675, 2676, 2677, 2678, 2679, 2680, 26

1	9	0	0	0
---	---	---	---	---

What would Isaac Newton do?



James Kwok

- 1 find descent direction
- 2 choose stepsize
- 3 descent

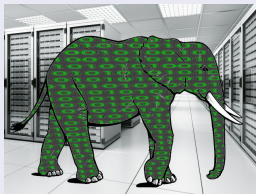
1000

- SVM: hinge loss (nonsmooth) + $\|w\|_2^2$ (smooth)
 - lasso: square loss (smooth) + $\|w\|_1$ (nonsmooth)
-
- extend gradient to nonsmooth functions \rightarrow subgradients

(Sub)Gradient Descent

Advantages

- easy to implement
- low per-iteration complexity → good scalability (BIG data)



Disadvantage

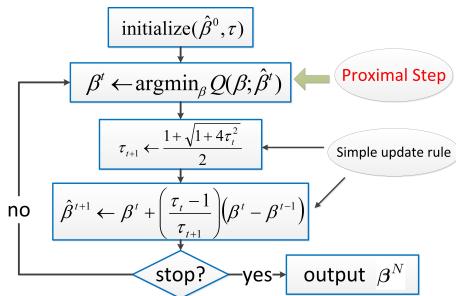
- uses **first-order** (gradient or subgradient) information
- slow convergence rate (especially for nonsmooth objectives)
→ may require a large number of iterations

FISTA (Beck & Teboulle, 2009)

Gradient descent on $f(w)$

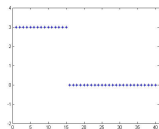
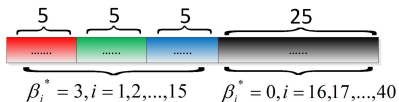
$$\hat{w}^t - \frac{1}{L} \nabla f(\hat{w}^t) = \arg \min_w (w - \hat{w}^t)^T \nabla f(\hat{w}^t) + \frac{L}{2} \|w - \hat{w}^t\|^2$$

For $f(w) + \Omega(w)$

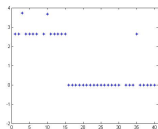


$$\arg \min_w Q(w, \hat{w}^t) \equiv (w - \hat{w}^t)^T \nabla f(\hat{w}^t) + \frac{L}{2} \|w - \hat{w}^t\|^2 + \Omega(w)$$

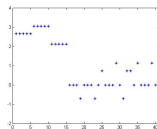
- L : Lipschitz constant of ∇f



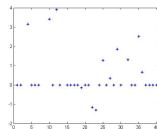
(a) ground truth



(b) OSCAR



(c) fused lasso



(d) lasso

- OSCAR can select relevant and correlated features

- (c) $\frac{1}{\sqrt{2}}$, $\frac{1}{\sqrt{2}}$, $\frac{1}{\sqrt{2}}$

		fused	elastic	
	lasso	lasso	net	OSCAR
test accuracy	65.9	72.4	71.8	74.1
#nonzero features	40.0	217.2	147.0	103.8

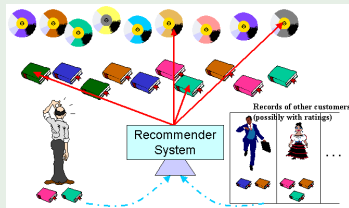
[Zhong & Kwok. Efficient sparse modeling with automatic feature grouping.
ICML-2011, and TNNLS-2012]

- often involves multiple learning tasks
- these tasks are **related** (share some information)

Example



digit recognition



product recommendation

- harness the task relationships \rightarrow learn all tasks **together**
- allow tasks to borrow strength from each other

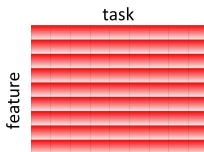
Table 1

$$\min_{w_1, \dots, w_T} \sum_{t=1}^T \left(\underbrace{\ell_t(w_t)}_{\text{loss}} + \underbrace{\Omega_t(w_t)}_{\text{task-specific regularizer}} \right) + \underbrace{\Omega_{\text{MTL}}(w_1, w_2, \dots, w_T)}_{\text{MTL regularizer}}$$

Different Assumptions for Ω_{MTL}

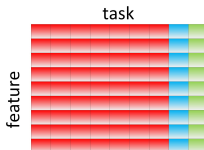
Assume: w_t 's form **one** group of correlated tasks

- regularized MTL



Assume: w_t 's form one cluster with a few outliers

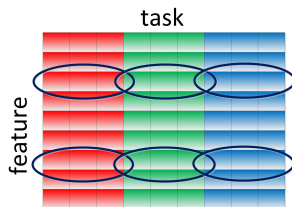
- robust MTL



Unlikely to hold with a **BIG** number of tasks

- clustered MTL

- clustered MTL



- needs to be fixed beforehand in clustered MTL

- needs to be fixed beforehand in clustered MTL

- flexible enough?

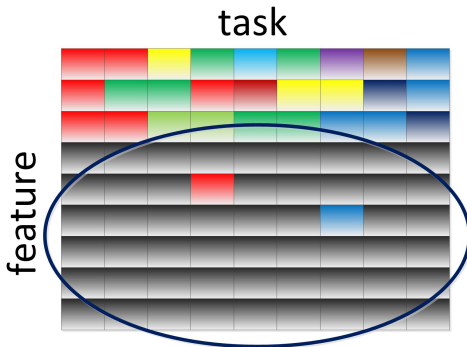
- flexible enough?

Motivating Example: Movie Recommendation



	A	B	C	D	E
Jackie Chan					
Action					
Comedy					
Kungfu					
Language					

Example 2: Features with Different Discriminating Power



Goal

Still want the tasks to be clustered, but

- task cluster structure can vary from feature to feature
- infer the clustering structure automatically

Idea

extend the feature clustering idea to task clustering in MTL

Flexible Task-Clustered MTL

- ① decompose each w_t into $u_t + v_t$
 - u_1, u_2, \dots, u_T : clustered
 - v_t : task-specific variation
- ② cluster the u_t 's feature by feature
 - for each feature d , minimize $|u_{i,d} - u_{j,d}|$ for all tasks i and j

$$\min_{U, V} \underbrace{\sum_{t=1}^T \|y_t - X_t(u_t + v_t)\|^2}_{\text{loss}} + \lambda_1 \|U\|_{\text{clus}} + \underbrace{\lambda_2 \|U\|_F^2 + \lambda_3 \|V\|_F^2}_{\text{ridge regularizers}}$$

- $U = [u_1, \dots, u_T], V = [v_1, \dots, v_T]$
- $\|U\|_{\text{clus}} = \sum_{d=1}^D \sum_{i < j} |u_{i,d} - u_{j,d}|$
 - a convex relaxation of **k-means** clustering on each feature

Special Cases

$$\min \sum_{t=1}^T \|y_t - X_t(u_t + v_t)\|^2 + \lambda_1 \|U\|_{clus} + \lambda_2 \|U\|_F^2 + \lambda_3 \|V\|_F^2$$

- ① $\lambda_1 = \lambda_2 = \lambda_3 = 0$: independent LS regression on each task
- ② $\lambda_1 = \infty$: regularized MTL (Evgeniou et al, 2005)
- ③ $\lambda_1 = 0$: independent ridge regression on each task

Nice Properties

Theoretical

With high probability, for large enough sample size,

- the obtained task weights converge to the ground truth
- U captures the clustering structure
 - tasks i, j are in the **same cluster** for feature $d \rightarrow u_{i,d} = u_{j,d}$
 - tasks i, j are in **different clusters** $\rightarrow u_{i,d} \neq u_{j,d}$

Computational

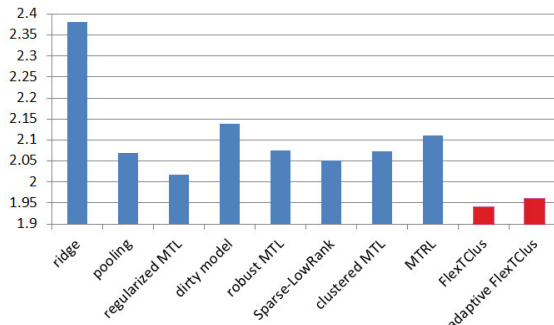
- optimization using **FISTA**
- similar to feature grouping, the **proximal step** can be efficiently computed in $O(TDn + DT \log T)$ time
 - T : number of tasks
 - D : feature dimension
 - n : sample size

Experiment: Product Rating

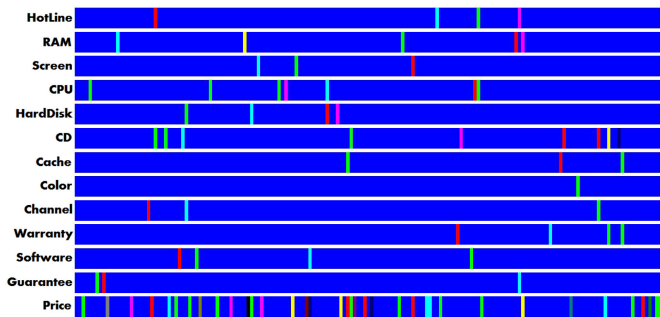


Predict the ratings of students (tasks) on personal computers (each described by 13 attributes)

Root mean squared error



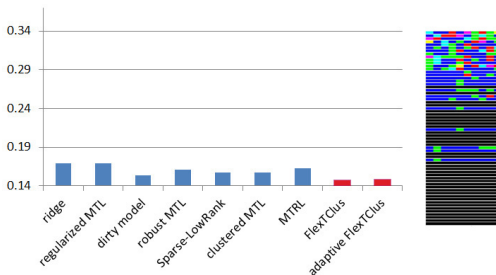
Product Rating: Task Clustering Structure



- one main cluster for the first 12 attributes (related to performance & service)
- lots of varied opinions on the last attribute (price)

Experiment: Digit Recognition

- 10-class classification problem \rightarrow 10 1-vs-rest problems
- use PCA features



- FlexTClus has the lowest classification error
- leading PCA features are more discriminative; trailing PCA features form one cluster close to zero (black)

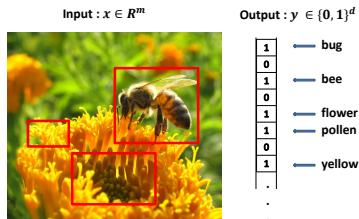
[Zhong & Kwok. [Convex multitask learning with flexible task clusters](#).

ICML-2012]

What if there are Many Many Tasks?

Multilabel classification

- an instance can have **more than one** labels



- cf. **multiclass** classification: an instance can have **only one** label
- basic approach: one label, one task (**binary relevance**)

Too Many Tasks

Flickr: > 20 million unique tags (labels) in 2010



animals architecture art asia australia autumn baby band barcelona beach berlin bike bird birds
 birthday black blackandwhite blue bw california canada canon car cat chicago
 china christmas church city clouds color concert dance day de dog england europe
 fall family fashion festival film florida flower flowers food football france friends
 fun garden geotagged germany girl graffiti green halloween hawaii holiday house india
 instagramapp iphone iphoneyography island italy japan kids
 lake landscape light live london love macro me mexico model museum music
 nature new newyork newyorkcity night nikon nyc ocean old paris park party
 people photo photography photos portrait raw red river rock san sanfrancisco
 scotland sea seattle show sky snow spain spring square squareformat
 street summer sun sunset taiwan texas thailand tokyo travel tree trees trip uk
 unitedstates urban usa vacation vintage washington water wedding white winter
 woman yellow zoo

Too Many Tasks...

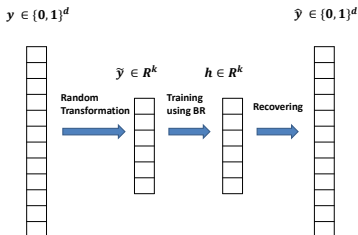
Open Directory Project: the largest human-edited Web directory

- over 4 million websites
- **787,774** categories (labels)

Top: Computers (94,925)

-
- | | |
|--|-------------------------------------|
| • Computer Science (1,876) | • Security (2,575) |
| • Hardware (5,205) | • Software (27,390) |
| • Internet (24,521) | • Systems (2,570) |
-
- | | |
|--|--|
| • Algorithms (303) | • Human-Computer Interaction (261) |
| • Artificial Intelligence (1,188) | • Intranet (39) |
| • Artificial Life (208) | • MIS@ (375) |
| • Bulletin Board Systems (93) | • Mobile Computing (486) |
| • CAD and CAM (731) | • Multimedia (2,455) |
| • Computer and Technology Law@ (105) | • Newsgroups@ (275) |
| • Data Communications (948) | • Open Source (699) |
| • Data Formats (1,679) | • Operating Systems@ (321) |
| • Desktop Publishing (113) | • Parallel Computing (321) |
| • E-Books (155) | • Performance and Capacity (43) |
| • Emulators (375) | • Programming (14,397) |
| • Fonts@ (272) | • Robotics (722) |
| • Games@ (31,921) | • Speech Technology (372) |
| • Graphics (1,243) | • Supercomputing (35) |
| • Hacking (167) | • Usenet (275) |
| • Home Automation (67) | • Virtual Reality (330) |

Label Transformation



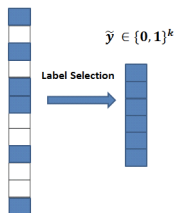
- ① **projects** the d -dimensional label vector to a k -dimensional space, where $k \ll d$
- ② learn a regression model for each dimension of the transformed label vector
- ③ predict in the low-dimensional space, then back-project to the d -dimensional space

The transformed labels, though fewer in quantity, may be more difficult to learn

Label Selection

- selects only a **few** output labels for training; reconstruct the other output labels from this label subset

$$\mathbf{y} \in \{0,1\}^d$$



- label subset comes from the original labels \rightarrow learning problems will not be more difficult

How to find the label subset?

- optimization (group-sparse learning problem)
- expensive, esp. with a lot of labels

Column Subset Selection Problem

Finding the label subset is a **column subset selection problem**

Column Subset Selection Problem (CSSP)

- given: matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$; positive integer k
- find C (k columns of \mathbf{A}) that spans \mathbf{A} as much as possible
 - \mathbf{A}_C : submatrix of \mathbf{A} with columns indexed by C

$$\min_{C: |C|=k} \|\mathbf{A} - \underbrace{\mathbf{A}_C \mathbf{A}_C^\dagger}_{\text{project onto } \mathbf{A}_C} \mathbf{A}\|_F$$

- in our case, \mathbf{A} is just the label matrix \mathbf{Y} !
 - \mathbf{Y} : each column is a label, each row is a sample

Existing CSSP Solvers

- perform **random sampling** of the columns
 - sampling probability is based on the **leverage** score

$$p_i = \|(\mathbf{V}_k^T)_{(i)}\|_2^2 / k$$
 - \mathbf{V}_k : top k right singular vectors of \mathbf{Y} ; $(\mathbf{V}_k^T)_{(i)}$: i th column of \mathbf{V}_k^T
 - leverage has been used to detect **outliers** in regression analysis
 → **importance** of each sample
- 1 (Drineas et al., 2006): sample $O(k^2)$ columns
 - a lot more than what we need (which is k)
 - 2 (Boutsidis et al., 2009)
 - sample $\Theta(k \log k)$ columns; then post-process to get k columns
 - post-processing can be even more computationally expensive than the sampling step itself!

Multilabel Classification via CSSP (ML-CSSP)

- 1: compute the sampling probability p_i for each column in \mathbf{Y} based on leverage;
- 2: $C \leftarrow \emptyset$;
- 3: **while** $|C| < k$ **do**
- 4: sample with replacement a column from \mathbf{Y} using p_i 's;
- 5: **if** $i \notin C$ **then**
- 6: $C \leftarrow C \cup \{i\}$;
- 7: **end if**
- 8: **end while**
- 9: train classifiers for the k selected labels.

Some Properties

With high probability, $\|\mathbf{Y} - \mathbf{Y}_C \mathbf{Y}_C^\dagger \mathbf{Y}\|_F \leq \text{constant} \times \|\mathbf{Y} - \mathbf{Y}_k\|_F$

- \mathbf{Y}_k : best rank- k approximation of \mathbf{Y}

With high probability, we can get k different columns in $O(k \log k)$ sampling rounds

- may still need to sample $O(k \log k)$ columns in the worst case
- empirical results show much smaller number

Time complexity

(Drineas et al., 2006)	$O(ndk)$	+	$O(k^2)$	
(Boutsidis et al., 2009)	$O(ndk)$	+	$O(k \log k)$	+
ours	$O(ndk)$	+	$O(k \log k)$	$O(k^3 \log^2 k \log(k \log k))$

Can also be [kernelized](#)

Experiments

Compare ML-CSSP with

- ① label transformation methods: PLST (Tai & Lin, 2012), CPLST (Chen & Lin, 2012) , CL (Zhou et al., 2012)
- ② label selection method: MOPLMS (Balasubramanian & Lebanon, 2012)
- ③ Binary relevance (BR)

data set	#samples	#features	#labels
cal500	502	68	174
corel5k	5,000	499	374
delicious	16,105	500	983
EUR-Lex (dc)	19,348	5,000	412
EUR-Lex (desc)	19,348	5,000	3,993
dmoz	394,756	829,208	35,437

RMSE and AUPRC

data set	RMSE					
	ML-CSSP	MOPLMS	PLST	CPLST	CL	BR
cal500	4.93	5.04	4.97	5.00	5.70	5.06
corel5k	1.89	1.89	1.91	1.91	2.71	1.91
delicious	4.29	-	4.27	4.26	5.58	4.26
EUR-Lex (dc)	1.22	-	1.22	1.23	2.03	1.50
EUR-Lex (desc)	2.93	-	3.02	3.06	4.51	3.51
dmoz	2.83	-	2.95	-	-	4.02

AUPRC (Area Under the Precision-Recall Curve)

data set	ML-CSSP	MOPLMS	PLST	CPLST	CL	BR
cal500	0.500	0.459	0.488	0.412	0.169	0.442
corel5k	0.089	0.080	0.079	0.082	0.011	0.083
delicious	0.220	-	0.182	0.227	0.089	0.237
EUR-Lex (dc)	0.180	-	0.180	0.167	0.036	0.173
EUR-Lex (desc)	0.094	-	0.018	0.086	0.016	0.086
dmoz	0.016	-	0.001	-	-	0.012

- many methods cannot be run because of the large number of labels (marked "-")
- ML-CSSP obtains better performance

Encoding Time

- time to perform label transformation/selection

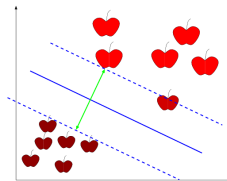
data set	ML-CSSP	MOPLMS	PLST	CPLST	CL
cal500	0.0	10.7	0.0	0.0	182.0
corel5k	0.2	36.8	0.2	0.3	292.0
delicious	11.6	-	11.6	16.0	5675.2
EUR-Lex (dc)	4.0	-	4.0	352.9	547.1
EUR-Lex (desc)	153.8	-	153.8	511.7	15585.5
dmoz	1428.9	-	1428.7	-	-

- ML-CSSP is almost as efficient as the fastest one (PLST)

[Bi & Kwok. [Efficient multi-label classification with many labels](#). ICML-2013]

Big Sample Size

ML tool: **Kernel methods**



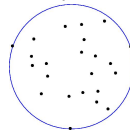
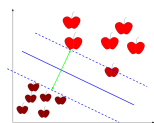
Example (big data problem)

SVM on n training examples:

- $O(n^2)$ memory for the $n \times n$ kernel matrix
- inverting/eigenvalue decomposition of the kernel matrix $\rightarrow O(n^3)$ time

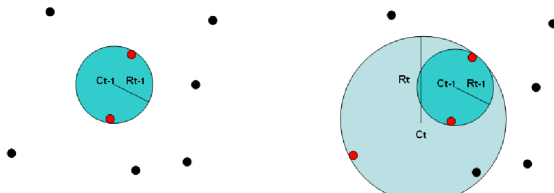
Core Vector Machine

- ① SVM training is a **minimum enclosing ball** (MEB) problem



- ② **near-optimal** solutions are good enough in practical applications → efficient **approximation algorithms**

- at each iteration, expand the current ball $B(c_t, r_t)$ by including the furthest point
- repeat until all the points are covered by $B(c_t, (1 + \epsilon)r_t)$



CVM in Classification

KDDCUP-99 intrusion detection data set (4,898,431 patterns)

method		# train patns input to SVM	# test errors	SVM training time (in sec)	other proc time (in sec)
random sampling	0.001%	47	25,713	0.000991	500.02
	0.01%	515	25,030	0.120689	502.59
	0.1%	4,917	25,531	6.944	504.54
	1%	49,204	25,700	604.54	509.19
	5%	245,364	25,587	15827.3	524.31
active learning		747	21,634	94192.213	
CB-SVM		4,090	20,938	7.639	4745.483
CVM		4,898,431	19,513	1.4	

- CVM is **fast** and **accurate**

[Tsang, Kwok, Cheung. [Core vector machines: Fast SVM training on very large data sets](#). **JMLR**, 2005]

Potential Limitations

Still need to solve a QP in finding the MEB of the core set

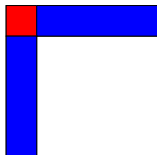
- very large data set \rightarrow large core set \rightarrow large QP



Low-Rank Approximation: Nyström Algorithm

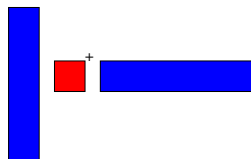
Sample m columns from matrix K

$$C = \begin{bmatrix} W \\ S \end{bmatrix} \quad K = \begin{bmatrix} W & S^T \\ S & B \end{bmatrix}$$



Rank- k Nyström approximation:

$$CW_k^+ C^T$$



Time complexity: $O(nmk + m^3)$

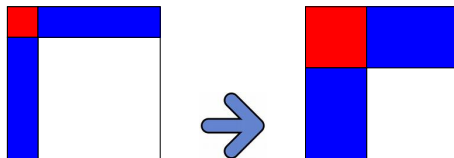
- $m \ll n \rightarrow$ much lower than the $O(n^3)$ complexity

How to Choose the Columns?

- ① random sampling
- ② probabilistic
 - chooses the columns based on a data-dependent probability
- ③ greedy approach
- ④ **clustering**-based
 - inexpensive; with interesting theoretical properties

[Zhang & Kwok. [Clustered Nystrom method for large scale manifold learning and dimension reduction](#). **TNN**, 2010]

Tradeoff between Accuracy and Efficiency



- more columns sampled, more accurate is the approximation

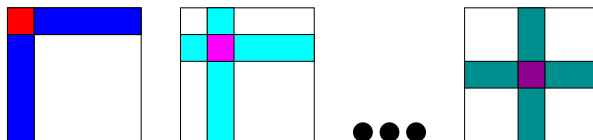
Example (data set with several millions samples)

- sampling only 1% of the columns
- W larger than $10,000 \times 10,000$
- SVD on W dominates and becomes prohibitive

Ensemble Nyström Algorithm (Kumar et al., 2009)

Replace the **large** SVD by a number of small SVDs

- **ensemble** of n_e Nyström approximators, each samples m columns



- each expert performs a standard Nyström approximation
- linearly **combine** the n_e rank- k approximations
 $\tilde{G}_{1,k}, \tilde{G}_{2,k}, \dots, \tilde{G}_{n_e,k}$

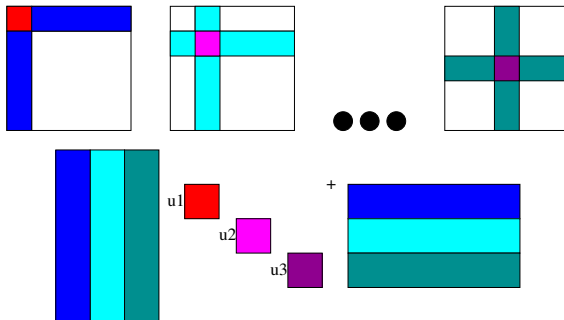
$$\tilde{G}^{ens} = \sum_{i=1}^{n_e} \mu_i \tilde{G}_{i,k} \quad (\mu_i \text{'s: mixture weights})$$

Ensemble Nyström Algorithm...

Time complexity

- $O(n_e n m k + n_e m^3 + C_\mu)$
 - C_μ : cost of computing the mixture weights
 - serial implementation: (roughly) n_e times that of Nyström
 - parallel implementation: (roughly) similar to Nyström

Implicitly, approximate $W^+ \in \mathbb{R}^{n_e m \times n_e m}$ by a **block diagonal** matrix

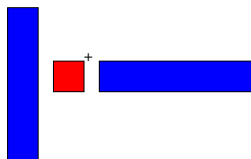


- less accurate

Nyström + Randomized SVD

Idea

- 1 sample **more columns** (like ensemble Nyström)
- 2 use a SVD approximation that is **more accurate** but still relatively efficient



- standard Nyström: SVD ($O(m^3)$ complexity)
- ensemble Nyström: block diagonal assumption
- our proposal: **randomized** SVD ($O(m^2k + k^3)$ complexity)

Is it Efficient?

Recall that $n \gg m \gg k$

Nyström



$$O(nmk + m^3)$$

ensemble Nyström



...



$$O(nmk + n_e k^3 + C_\mu)$$

randomized SVD

$$O(n^2 k + k^3)$$

Nyström + rand. SVD

$$O(nmk + m^2 k + k^3) = O(nmk + k^3)$$



Is it Accurate?

Rank- k standard Nyström approximation \hat{K} (with m randomly sampled columns)

$$\mathbb{E} \|K - \hat{K}\|_2 \leq \|K - K_k\|_2 + \frac{2n}{\sqrt{m}} (\max_i K_{ii})$$

- K_k : best rank- k approximation

Proposed method

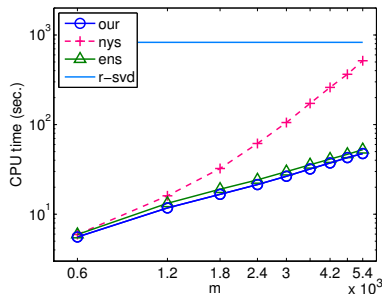
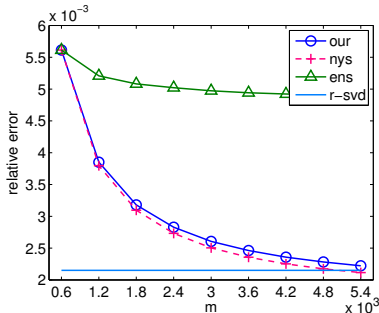
$$\mathbb{E} \|K - \hat{K}\|_2 \leq \zeta^{1/q} \|K - K_k\|_2 + (1 + \zeta^{1/q}) \frac{n}{\sqrt{m}} (\max_i K_{ii})$$

- ζ : constant and $\zeta^{1/q}$ close to 1
 \rightarrow becomes $\|K - K_k\|_2 + \frac{2n}{\sqrt{m}} (\max_i K_{ii})$, same as that for standard Nyström

Proposed method is as accurate as standard Nyström

Experiment

x-axis: number of sampled columns



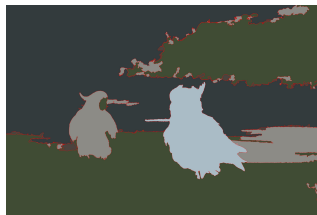
	accuracy	cost
randomized SVD	✓	✗
standard Nyström	✓	✗
ensemble Nyström	✗	✓
proposed method	✓	✓

Image Segmentation

- Intel Xeon X5540 CPU, 16GB memory (matlab)



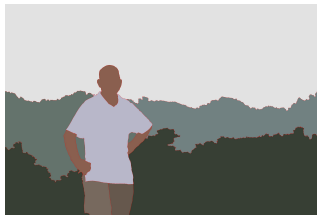
3872×2592 (10M pixels)



16.8 sec



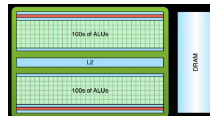
4752×3168 (15M pixels)



22.6 sec

Graphics Processors (GPU)

Popularly used in entertainment, high-performance computing



NVIDIA Tesla C1060

- 240 streaming processor cores
- peak single-precision (SP) performance: 933 GFLOPS
- peak double-precision (DP) performance: 78 GFLOPS

Intel Core i7-980X CPU

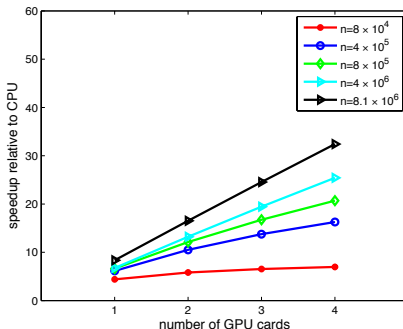
- 6 cores
- peak SP: 158.4 GFLOPS
- peak DP: 79.2 GFLOPS

Experiments on GPU

Machine

- two Intel Xeon X5560 2.8GHz CPUs, 32G RAM
- four NVIDIA Tesla C1060 GPU cards

Linear speedup with number of GPUs (MNIST-8M data set)



Different Numbers of Sampled Columns m

m	time (sec)			relative approx. error
	run on CPU	run on 4 GPUs	(speedup)	
2,000	908.1	24.5	(37.1x)	5.9×10^{-4}
4,000	1,789.4	47.2	(37.9x)	5.5×10^{-4}
6,000	2,647.8	81.7	(32.4x)	5.3×10^{-5}
8,000	3,556.5	104.8	(33.9x)	5.4×10^{-5}
10,000	4,426.4	119.5	(37.0x)	3.0×10^{-5}
20,000	8,988.2	253.8	(35.4x)	1.1×10^{-5}

More GPU Comparisons

- GPU versions of standard Nyström and ensemble Nyström

Time

number of samples (n)	Nyström	ours	ensemble Nyström
4×10^5	244.9	12.7	25.3
8×10^5	368.8	22.4	39.3
4×10^6	2,041.7	94.2	46.5
8.1×10^6	6,884.0	81.7	81.2

Relative approximation error

n	Nyström	ours	ensemble Nyström
4×10^5	8.5×10^{-5}	9.4×10^{-5}	4.4×10^{-4}
8×10^5	8.2×10^{-5}	8.3×10^{-5}	4.5×10^{-4}
4×10^6	9.0×10^{-6}	9.1×10^{-6}	6.9×10^{-5}
8.1×10^6	7.8×10^{-6}	7.8×10^{-6}	6.9×10^{-5}

	accuracy	cost
standard Nyström	✓	✗
ensemble Nyström	✗	✓
proposed method	✓	✓

Conclusion

Big data is **difficult** to handle

- high-dimensional
- involves a lot of tasks (related in some complex manner)
- big sample size
- **high velocity**

But there is **hope**

- better **models**
 - OSCAR, flexible task clustering, **transfer learning**
- better **tools**: optimization solvers and approximations
 - accelerated gradient descent, column subset selection solvers, Nyström algorithm, **stochastic algorithms**, **online algorithms**
- better **hardware**
 - GPU, **parallel/distributed architecture**

Isaac Newton

Modern day Isaac Newton