



CLASSICAL RANDOM FEATURE HASHING [GANCHEV AND DREDZE, 2008, SHI ET AL., 2009, WEINBERGER ET AL., 2009]

Consider **classification task** (for simplicity):

→ sparse **high-dimensional** labeled data:
 $(\phi_n, y_n) \in \mathbb{R}^D \times \{-1, +1\}, D \gg 1$

→ **linear** scoring $f(\phi; \mathbf{w}) = \langle \mathbf{w}, \phi \rangle$

✓ high dimension $D \Rightarrow$ approximate separability

✗ **But we need to access \mathbf{w} rapidly:**

→ must keep \mathbf{w} in RAM \rightarrow may not fit

→ storage model matters:

- (un)ordered associative tables need RAM and/or are slow

✓ **linear arrays are fastest** \Rightarrow
need integer feature indexes

Common solution:

• **feature hashing trick** (alphabet elimination/random feature mixing)

• **Idea:** – per-coordinate mapping into a lower-dimension feature space
– with **data-independent** pseudo-random function $HASH$:

$$\phi'_{d'} = \sum_{d: HASH(d)=d'} \phi_d, \quad \text{where } d \text{ is a feature key (usually a string)}$$

✓ $\|\phi'\| \simeq \|\phi\|$ with high probability

✓ works surprisingly well: **little or no sacrifice in quality!**

✓ implemented in many learning kits (Vowpal Wabbit, etc.)

Can we do even better with **data-dependant** hashing ?

PAPER HIGHLIGHTS

✓ replaces random feature hashes with **learned feature hashes**

✓ **simple greedy algorithm** for learning hashes

✓ learning **optimizes task objective**

✓ data-dependent hashing **improves classification** accuracy:

★ for high-dimensional features

★ in case of tight memory constrains (many collisions)

New hash function $HASH_{new}$ that:

• is informed of the final learning task

• can be done in preprocessing

• leverages existing optimization procedures

IDEA OVERVIEW

Intuition: hashing d, d' together entails less loss change if $w_d \simeq w_{d'}$
 \Rightarrow approximate $w_d - w_{d'}$ as dist. between learnable representations:

0 Equip each feature d with a vector $\nu(d) \in \mathbb{R}^V$
(e.g. some NLP stats about d 's usage in the wild)

1 Learn a map $\nu(d) \mapsto \mathcal{H}(\nu(d)) \in \{0, 1\}^T$, s.t. the Hamming dist. between $\mathcal{H}(\nu(d_1)), \mathcal{H}(\nu(d_2))$ captures d_1, d_2 's similarity for task

2 Apply surjection $\mathcal{B} : \{0, 1\}^T \rightarrow \{0 \dots M\}$, s.t. close Hamming vectors get projected into the same integer

3 Define $HASH_{new} = \mathcal{H} \circ \mathcal{B}$, interpret outputs as RAM addresses.

LEARNING FEATURE HASHES BY OPTIMIZING HINGE LOSS WITH BOOSTING

1: greedy learning of Hamming representation \mathcal{H}

- will look for w_d representable as $\sum_{t \leq T} \alpha_t h_t(\nu(d))$
- **close** $\mathcal{H}(\nu(d)) = [h_1(\nu(d)), \dots, h_T(\nu(d))]$ for different $d \Rightarrow$ **close** values of respective w_d .

Hinge loss: $L = \sum_n (1 - y_n \sum_d \sum_{t \leq T} \alpha_t h_t(\nu(d)) \phi_{n,d})_+$

Learning $\mathcal{H}(\nu_d)$ in a **boosting fashion:**

- $\mathcal{H}^{t-1}(\nu) = [h_1(\nu), \dots, h_{t-1}(\nu)] \in \{0, 1\}^{t-1}$
- $\mathcal{H}^t(\nu)$ is obtained by appending a bit-function – a per-coordinate decision stump $h_t = \llbracket \nu_{k^*} > \theta^* \rrbracket$:

$$k^*, \theta^* = \arg \max_{k, \theta} \left| \sum_{n: y_n \langle \mathbf{w}^t, \phi_n \rangle < 1} y_n \sum_d \phi_{n,d} \llbracket \nu_{d,k} \geq \theta \rrbracket \right|$$

2: distance-sensitive projection \mathcal{B}

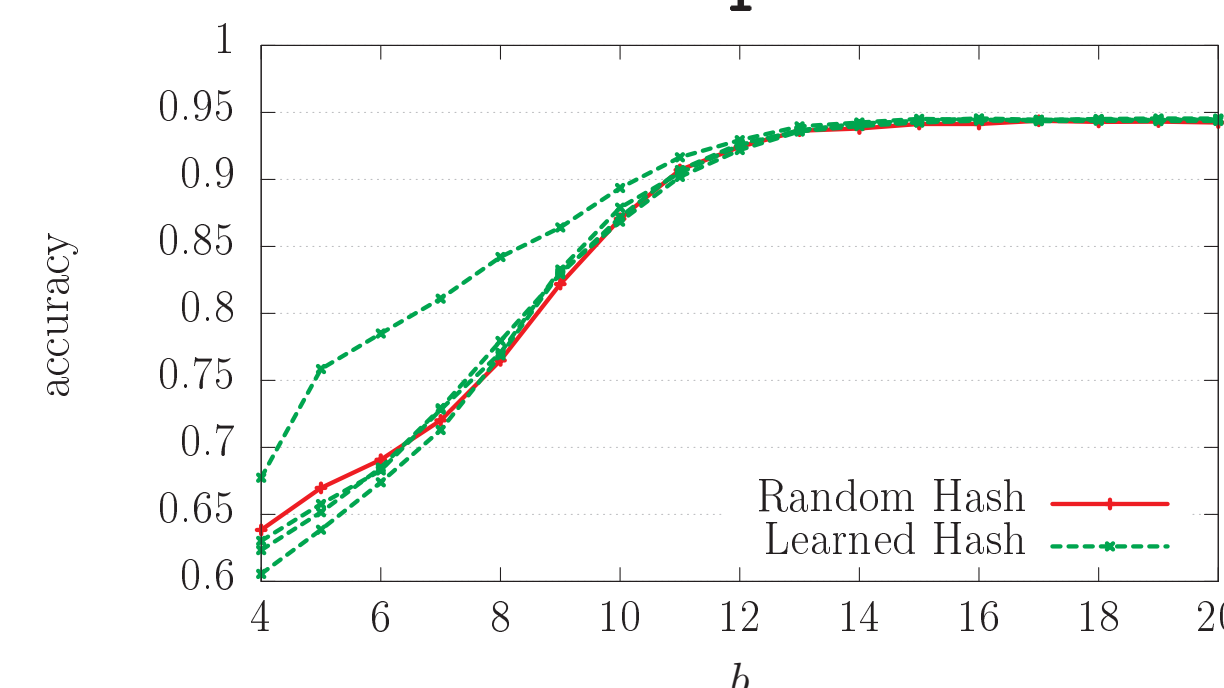
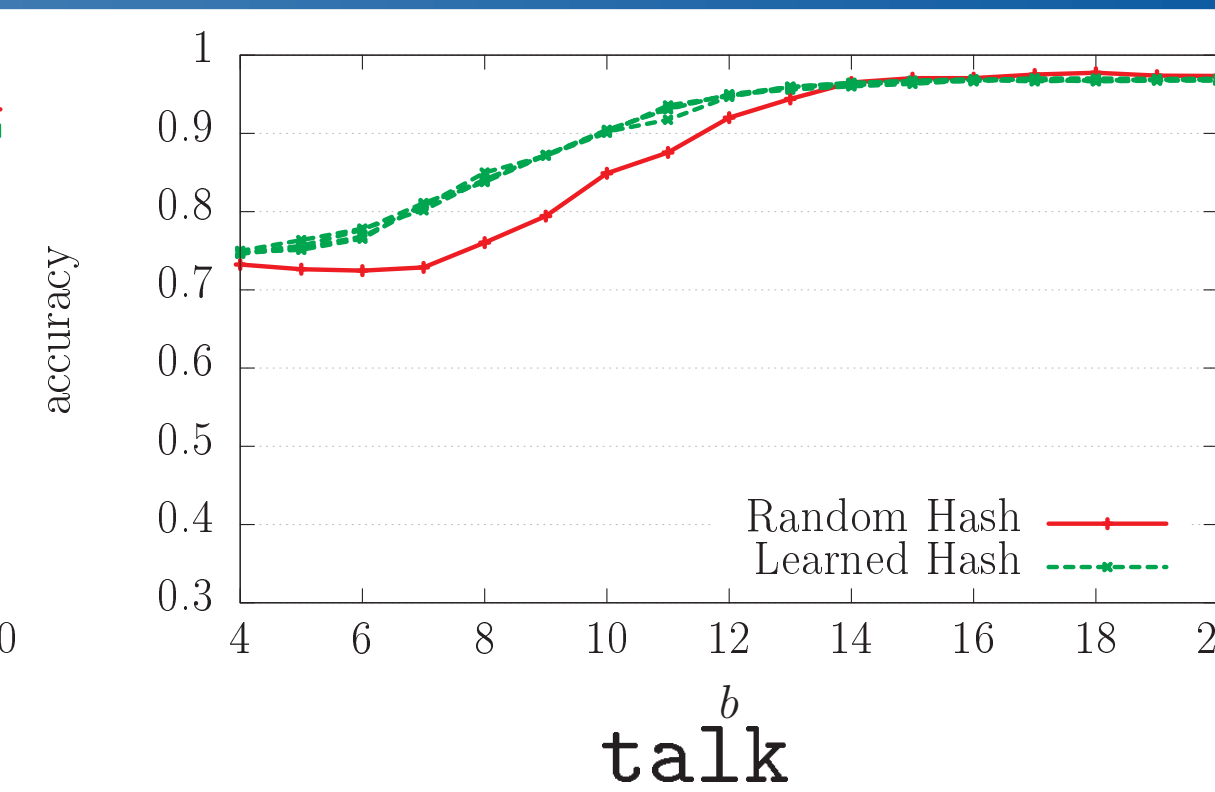
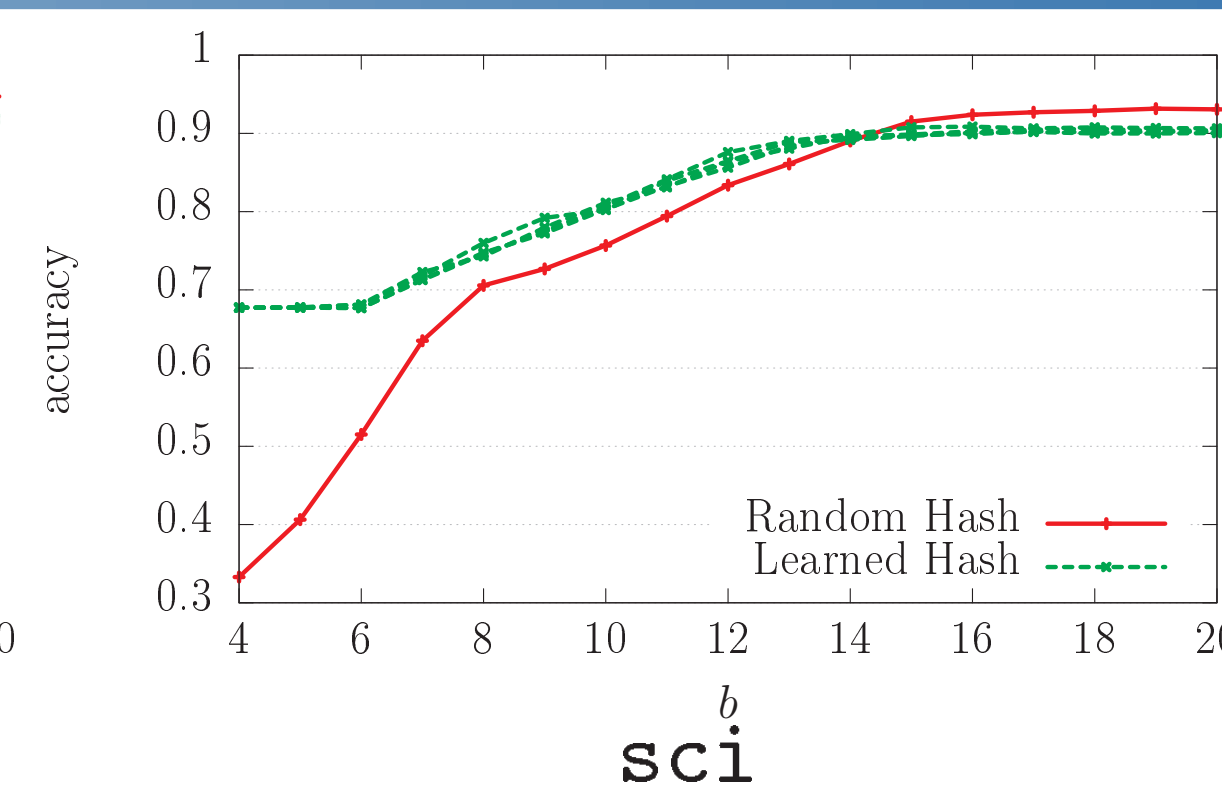
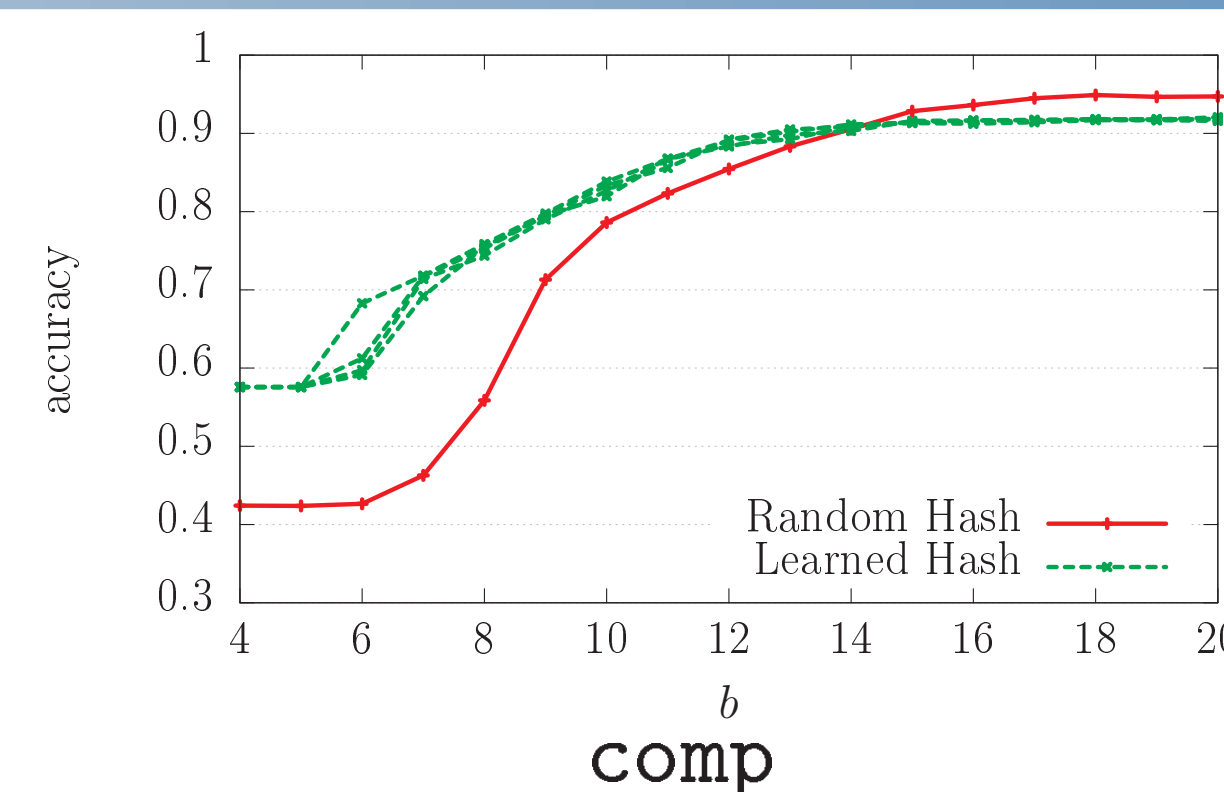
- **Task:** compress $\mathcal{H}(\nu)$ into short codes that have a high collision probability for close $\mathcal{H}(\nu)$.
- we use the **KOR random traces** [Kushilevitz et al., 1998]:
 - for a bit-vector $\mathbf{h} = [h_1, \dots, h_T]$
 - and random Bernoulli vectors $\mathbf{r}_m = [r_{m,1}, \dots, r_{m,T}]$
 - the trace is $\mathbf{t} = [t_1, \dots, t_M]$, s.t. $t_m = \langle \mathbf{h}, \mathbf{r}_m \rangle \bmod 2$.
- t_m has a bias towards closer \mathbf{h} , amplified by repeating M times
- **collision probability decays** with $D_H(\mathbf{h}_1, \mathbf{h}_2)$ and M :

$$P[t_1 = t_2 | D_H(\mathbf{h}_1, \mathbf{h}_2) \leq \Delta] \geq \left(\frac{1}{2} + \frac{1}{2}(1 - 2p)^\Delta\right)^M$$

PROOF-OF-CONCEPT EXPERIMENTS

★ dataset #1: **20-newsgroups**

- 3 (one vs. all) class. tasks for comp, sci, talk
- 70% train / 30% test
- feature dim.: **700K** (all 2-grams)



★ dataset #2: **RCV1**

- hashes learned on 100K examples
- Vowpal Wabbit ran on full train set
- feature dim.: **40K**

- vectors ν filled with 1/2-gram **DICE** coefficients: $\nu'_d(d) = \text{dice}(d, d')$ for all d, d'
- Vowpal Wabbit was trained on preprocessed data (with learned hashes)
- plots: classification **accuracy vs. # of bits b**
- **the higher the curve, the better**

CONCLUSION

For the considered **classification** task:

- the less RAM is available (the smaller b)
- or the higher dim. input feature space has

.. **the more sense it makes to learn hashes**

Future work:

- ranking objectives
- tasks with very large dimension:
 - information retrieval
 - collaborative filtering

References

- [Ganchev and Dredze, 2008] Ganchev, K. and Dredze, M. (2008). Small statistical models by random feature mixing. In *Proceedings of the ACL-2008 Workshop on Mobile Language Processing*. Association for Computational Linguistics.
- [Kushilevitz et al., 1998] Kushilevitz, E., Ostrovsky, R., and Rabani, Y. (1998). Efficient search for approximate nearest neighbor in high dimensional spaces. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, STOC ’98, pages 614–623, New York, NY, USA. ACM.
- [Shi et al., 2009] Shi, Q., Petterson, J., Dror, G., Langford, J., Smola, A. J., and Strehl, A. L. (2009). Hash Kernels. 5:496–503.
- [Weinberger et al., 2009] Weinberger, K., Dasgupta, A., Langford, J., Smola, A. J., and Attenberg, J. (2009). Feature hashing for large scale multitask learning. *Computing Research Repository*, abs/0902.2:140–1120.