# Fully Understanding the Hashing Trick

**Lior Kamma, Aarhus University**

Joint work with **Casper Freksen** and **Kasper Green Larsen.**

**AARHUS UNIVERSITY**

# Recommendation and Classification

PG-13

Comic Book

Super Hero

Sci Fi

Adventure

Action

Violent

Scary

Comedy

Drama

Horror

# Recommendation and Classification

PG-13

Comic Book

Super Hero

Sci Fi

Adventure

Action
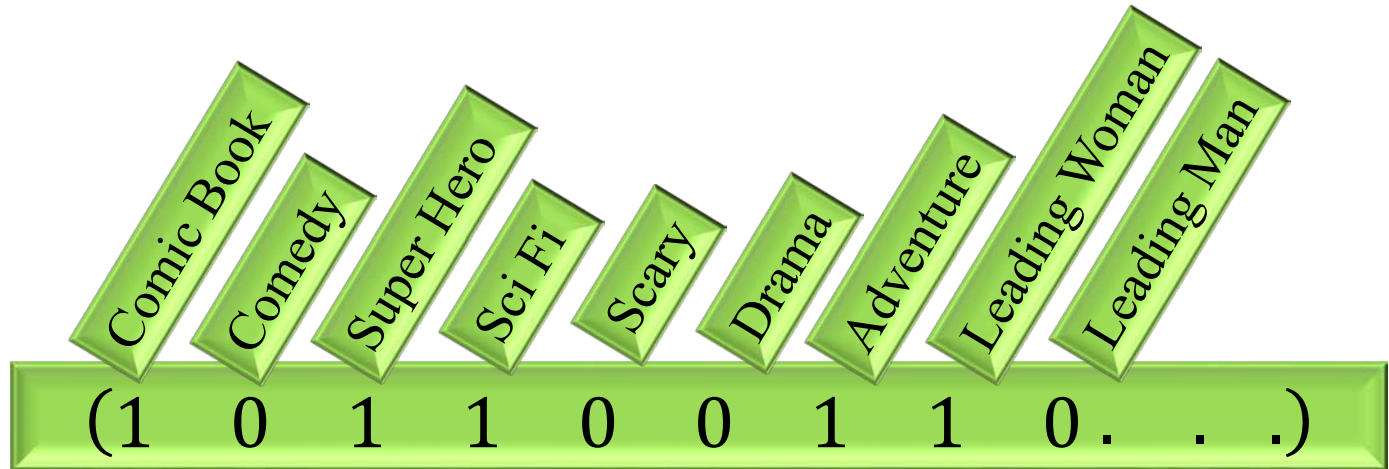
Violent

Scary

Comedy

Drama

Horror

Categorical Variables

How do we decide these are "close"?

# Feature Vectors

Comic Book    Comedy    Super Hero    Sci Fi    Scary    Drama    Adventure    Leading Woman    Leading Man

$$(1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad 1 \quad 0. \quad . \quad .)$$

$$(1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad 0 \quad 1. \quad . \quad .)$$

Boolean vectors

Denote the *feature dimension* by $n$

# $k$-Nearest Neighbours

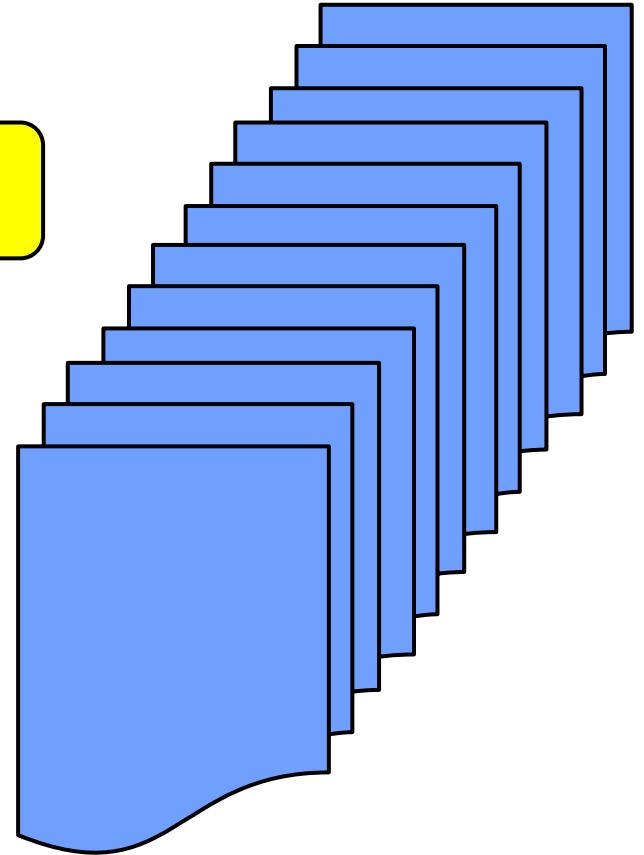Storing a corpus of $M$ items requires $\Omega(nM)$ memory

Corpus

# $k$-Nearest Neighbours

New Movie

**How do we find the $k$ closest movies?**

# Dimensionality Reduction

■ Given $\varepsilon, \delta \in (0,1)$ find

Approximation Ratio

Error Probability

# Dimensionality Reduction

Given $\varepsilon, \delta \in (0, 1)$ find random $f: \mathbb{R}^n \to \mathbb{R}^m$ such that for every $x, y \in \mathbb{R}^n$

For some small $m$

Think of $n$ as **HUGE**

# Dimensionality Reduction

- Given $\varepsilon, \delta \in (0,1)$ find random $f: \mathbb{R}^n \to \mathbb{R}^m$ such that for every $x, y \in \mathbb{R}^n$

$$\Pr[\|f(x) - f(y)\|_2^2 \in (1 \pm \varepsilon)\|x - y\|_2^2] \geq 1 - \delta$$

# Dimensionality Reduction

- Given $\varepsilon, \delta \in (0,1)$ find random $A \in \mathbb{R}^{m \times n}$ such that for every $x, y \in \mathbb{R}^n$

$$\Pr[\|A(x - y)\|_2^2 \in (1 \pm \varepsilon)\|x - y\|_2^2] \geq 1 - \delta$$

Focus on linear projections

Why linear?

- Cool Math

- Streaming (updates).

- Good in practice

UNIVERSITY

# Dimensionality Reduction

- Given $\varepsilon, \delta \in (0,1)$ find random $A \in \mathbb{R}^{m \times n}$ such that for every $x \in \mathbb{R}^n$

$$\Pr[\|A(x)\|_2^2 \in (1 \pm \varepsilon)\|x\|_2^2] \geq 1 - \delta$$

Focus on linear projections

Why linear?

- Cool Math

- Streaming (updates).

- Good in practice

UNIVERSITY

# Johnson Lindenstrauss Lemma [JL'84]

- Given $\varepsilon, \delta \in (0,1)$ there exists a random linear $A \in \mathbb{R}^{m \times n}$ such that for every $x$

$$\Pr[\|A(x)\|_2^2 \in (1 \pm \varepsilon)\|x\|_2^2] \geq 1 - \delta$$

$$m = O\left(\frac{\lg 1/\delta}{\varepsilon^2}\right)$$

In most proofs matrix is as dense as possible. Embedding takes $O(mn)$ operations.

AARHUS UNIVERSITY

# Johnson Lindenstrauss Lemma [JL'84]

- Given $\varepsilon, \delta \in (0,1)$ there exists a random linear $A \in \mathbb{R}^{m \times n}$ such that for every $x$

$$\Pr[\|A(x)\|_2^2 \in (1 \pm \varepsilon)\|x\|_2^2] \geq 1 - \delta$$

If $A$ is sparse, this can be made faster.

In most proofs matrix is as dense as possible. Embedding takes $O(mn)$ operations.

AARHUS UNIVERSITY

# Feature Hashing [Weinberger *et al.* 2009]

Add random signs

General Idea: Shuffle the entries of $x$

$x$

$$(1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad 0 \quad 1)$$

**+**

**-**

# Feature Hashing [Weinberger *et al.* 2009]

General Idea: Shuffle the entries of $x$

Add random signs

$x$     (1   0   1   1   0   0   1   0   1)

**+**

**-**

$f(x)$    0    1    0

$m = 3$

# Feature Hashing [Weinberger *et al.* 2009]

Add random signs

General Idea: Shuffle the entries of $x$

$x$     $(1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad 0 \quad 1)$

$+$

$-$

$f(x)$    0    1    $-1$    $m = 3$

# Feature Hashing [Weinberger *et al.* 2009]

General Idea: Shuffle the entries of $x$

Add random signs

**Observation: This operation is linear.**

**Moreover, every column has exactly one non-zero entry.**

$f(x)$

$0$   $1$   $-1$

$m = 3$

# The Hashing Trick – With High Prob.

- Observation: If $m$ is large enough, and the "mass" of x is not concentrated in few entries, then the trick works with high probability.

$$\varepsilon = 0.1$$

$$\begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\Pr_{h:\{1,2,\ldots,n\}\rightarrow\{1,2,\ldots,m\}}[h(1) = h(2)] = \frac{1}{m}$$

$$\frac{\|x\|_\infty}{\|x\|_2} = \frac{1}{\sqrt{2}}.$$

# The Hashing Trick – With High Prob.

Success iff no collision occurs nough, and the
m̶a̶s̶s̶ x is not concentrated in few entries,
the̶ trick works with high probability.

$\varepsilon = 0.1$

$$\begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\Pr_{h:\{1,2,\dots,n\}\to\{1,2,\dots,m\}}[h(1) = h(2)] = \frac{1}{m}$$

$$\frac{\|x\|_\infty}{\|x\|_2} = \frac{1}{\ }$$

To succeed we need $m \geq \frac{1}{\delta}$

# Tight Bounds – Formal Problem

- Fix $m, \varepsilon, \delta$.

- Define $\nu(m, \varepsilon, \delta)$ to be the maximum $\nu$ such that whenever $\|x\|_\infty \leq \nu \|x\|_2$ then feature hashing works.

# Tight Bounds – Formal Problem

- Fix $m, \varepsilon, \delta$.

- Define $\nu(m, \varepsilon, \delta)$ to be the maximum $\nu$ such that whenever $\|x\|_\infty \leq \nu \|x\|_2$ then feature hashing

We have a fixed budget, and a fixed room for error.

Evaluating $\nu$ has been an open question for almost a decade.

# Tight Bounds – Our Result

- Fix $m, \varepsilon, \delta$.

**Theorem.**

1. If $m < \dfrac{c \ \log\frac{1}{\delta}}{\varepsilon^2}$ then $\nu = 0$.

Essentially, this means our budget is too small to do anything meaningful.

# Tight Bounds – Our Result

- Fix $m, \varepsilon, \delta$.

**Theorem.**

1. If $m < \dfrac{c \, \log\frac{1}{\delta}}{\varepsilon^2}$ then $\nu = 0$

2. If $m \geq \dfrac{2}{\delta\varepsilon^2}$ then $\nu = 1$ .

> Essentially, this means our budget is rich enough to do anything.

# Tight Bounds – Our Result

- Fix $m, \varepsilon, \delta$.

**Theorem**

> This is tight,
> which means this is the *right*
> expression.

$$\frac{\log\frac{1}{\delta}}{\varepsilon^2} \le m < \frac{1}{\delta\varepsilon^2} \quad \text{then}$$

$$\nu = \Theta\left(\sqrt{\varepsilon} \cdot \min\left\{ \frac{\log\dfrac{\varepsilon m}{\log\frac{1}{\delta}}}{\log\frac{1}{\delta}}, \sqrt{\frac{\log\dfrac{\varepsilon^2 m}{\log\frac{1}{\delta}}}{\log\frac{1}{\delta}}} \right\}\right)$$

# Empirical Analysis

Results show that the Θ-constant is close to 1.

This implies that Feature Hashing's performance can be very well predicted in practice using our formula.

$$\frac{\nu}{\sqrt{\varepsilon} \ \min\left\{\left(\dfrac{\lg \frac{\varepsilon m}{\lg 1/\delta}}{\lg 1/\delta}\right), \sqrt{\dfrac{\lg \frac{\varepsilon^2 m}{\lg 1/\delta}}{\lg 1/\delta}}\right\}}$$

$$\nu = \Theta\left(\sqrt{\varepsilon} \cdot \min\left\{\frac{\log \frac{\varepsilon m}{\log \frac{1}{\delta}}}{\log \frac{1}{\delta}}, \sqrt{\frac{\log \frac{\varepsilon^2 m}{\log \frac{1}{\delta}}}{\log \frac{1}{\delta}}}\right\}\right)$$

0.725

AARHUS UNIVERSITY

# Questions?



Come see poster

Read the paper

Talk offline

All of the above

Tight Cell-Probe Bounds for Succinct Boolean Matrix-Vector Multiplication

# Questions?

Come see poster

Read the paper

Talk offline

All of the above

Thank you