

# The unreasonable effectiveness of feature hashing

---

Gianluca Campanella [dʒanˈluːka kampaˈnella]

14<sup>th</sup> July 2019

I'm a Data Scientist at



**Microsoft**

in **AzureCAT**

## What I do nowadays

I also run my own company



Estimand

that provides

**Data Science consulting and training**

# Contents

---

Background

Feature hashing

Library support

# Acknowledgements

- Sharat Chikkerur
- Markus Cozowicz

# Background

---

# Supervised ML

- Observe  $\mathbf{X}_{n \times p}$  and  $\mathbf{y}_n$  (typically  $n \gg p$ )
- Find the ‘best’ estimator of some statistic of  $\mathbf{y}$  given  $\mathbf{X}$  (for example  $\mathbb{E}[\mathbf{y}|\mathbf{X}]$ )

# Supervised ML

- Observe  $\mathbf{X}_{n \times p}$  and  $\mathbf{y}_n$  (typically  $n \gg p$ )
- Find the ‘best’ estimator of some statistic of  $\mathbf{y}$  given  $\mathbf{X}$  (for example  $\mathbb{E}[\mathbf{y}|\mathbf{X}]$ )

## Generalised linear models (GLMs)

- $\mathbb{E}[\mathbf{y}|\mathbf{X}] = g^{-1}(\mathbf{X}\beta)$  with  $g$  given
- Find the ‘best’ estimates for  $\beta$



# Features are often categorical

- Discretised continuous features
- One-hot encoding
- Bag-of-words representation

# One-hot encoding

country	ENG	NIR	SCT	WLS
England	1	0	0	0
Wales	0	0	0	1
Scotland	0	0	1	0
Scotland	0	0	1	0
Northern Ireland	0	1	0	0

# Bag-of-words representation

---

document

---

Ashley loves cats. She also likes dogs.  
Barbara hates cats but she loves dogs.  
Carol loves cats and dogs.

---



	cat	dog	hate	like	love
1	1	1	0	1	1
2	1	1	1	0	1
3	1	1	0	0	1

# High-dimensional feature space

What's going on?

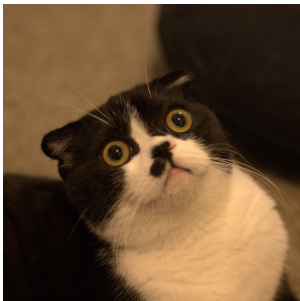
- $n$  is fixed
- $p$  depends on the cardinality of features

# High-dimensional feature space

## Why is this a problem?

- Online learning
- $p \rightarrow n$  breaks GLMs (though regularisation helps)

# High-dimensional feature space



cat      dog  
 $\begin{bmatrix} 0.99 & 0.01 \end{bmatrix}$

# High-dimensional feature space

Observations in  $X$  are high-dimensional but **sparse**

- Ad tech
- E-commerce
- Social networks

# High-dimensional feature space

## Can we find a 'smaller' $X$ ?

- Time- and space-efficient to compute
- Fixed  $p$
- Sparse



- Categorical and textual features are common
  - $p \propto$  feature cardinality (often very large)
  - Sparse  $\mathbf{X}$
- Want an efficiently computed, fixed- $p$ , sparse representation

# Feature hashing

---

# Feature hashing

country = 'England'



Hash value

$h(\text{country}=\text{England}) = 1462978411$



Feature index

$i = h(\text{country}=\text{England}) \bmod 100 = 11$

# Feature hashing

country	$i$	$\dots$	$x_8$	$\dots$	$x_{11}$	$\dots$	$x_{54}$	$\dots$	$x_{62}$	$\dots$
England	11				1					
Wales	62								1	
Scotland	54						1			
Scotland	54						1			
Northern Ireland	8		1							

# Feature hashing

- Time- (%timeit: 468 ns  $\pm$  7.2 ns) and space-efficient
- Fixed  $p$
- Sparse

# Hash function

$h$  maps **any** integer input onto integers in some range (e.g. `int32`)

## Important properties

- Uniform output
- Avalanche effect
  - $h(a) \rightarrow 3001393763$
  - $h(c) \rightarrow 1701913768$

# Feature hashing in Python

```
1 import mmh3
2
3 s = "country=England"
4
5 h = mmh3.hash(s, seed=42, signed=False)
6 print(h)    # => 1462978411
7
8 i = h % 100
9 print(i)    # => 11
```

# Hashing of Unicode strings

- Strings are just (very large) integers
- Be careful when handling Unicode

```
1 s1 = "sch\u00f6n"
2 s2 = "scho\u0308n"
3
4 print(s1)          # => schön
5 print(s2)          # => schön
6 print(s1 == s2)    # => False
```



# Unicode normalisation

```
1 import unicodedata
2
3 s1 = unicodedata.normalize("NFKD", "sch\u00f6n")
4 s2 = unicodedata.normalize("NFKD", "scho\u0308n")
5
6 print(s1)          # => schön
7 print(s2)          # => schön
8 print(s1 == s2)    # => True
```

# Projection

- The modulo operator is expensive
- If the hash size is a power of two, we can do better

```
1 %timeit h % 128
2 # => 81.3 ns ± 1.92 ns per loop
3
4 %timeit h & 127
5 # => 56.9 ns ± 1.4 ns per loop
```

# Collisions

```
1 def h(x):  
2     return mmh3.hash(x, seed=42, signed=False) & 127  
3  
4 print(h("country=United Kingdom")) # => 6  
5 print(h("country=Bulgaria"))      # => 6
```

- Smaller hash size → more collisions
- Impact on statistical performance and interpretability

# Sign function $\xi$

- Use another function  $\xi$  to determine the sign
- Collisions cancel out (in expectation)

```
1 def h(x):  
2     hash_ = mmh3.hash(x, seed=42, signed=True)  
3     return abs(hash_) & 127, (hash_ >= 0) * 2 - 1  
4  
5 print(h("country=United Kingdom")) # => ( 6, 1)  
6 print(h("country=Bulgaria"))       # => (122, -1)
```

# Recap

- $h$  maps any integer input onto integers in some range
- $h$  deterministically scrambles arbitrary-length bitmaps, producing fixed-length bitmaps
- The impact of collisions on statistical performance is small

# Feature hashing

---

## Example

## Original data

user	movie	has_cats	has_dogs	rating
Ashley	A Dog's Journey	0	1	4
Ashley	Catwoman	1	0	5
Ashley	The Aristocats	1	1	5
Ashley	The Queen's Corgi	0	1	3
Barbara	A Dog's Journey	0	1	5
Barbara	The Aristocats	1	1	2
Barbara	The Queen's Corgi	0	1	5
Carol	Catwoman	1	0	5
Carol	The Aristocats	1	1	5
Carol	The Queen's Corgi	0	1	4

## After feature hashing

user	movie	has_cats	has_dogs
2 <sup>-</sup>	5 <sup>-</sup>		110 <sup>-</sup>
2 <sup>-</sup>	1 <sup>-</sup>	9 <sup>+</sup>	
2 <sup>-</sup>	104 <sup>-</sup>	9 <sup>+</sup>	110 <sup>-</sup>
2 <sup>-</sup>	67 <sup>+</sup>		110 <sup>-</sup>
19 <sup>+</sup>	5 <sup>-</sup>		110 <sup>-</sup>
19 <sup>+</sup>	104 <sup>-</sup>	9 <sup>+</sup>	110 <sup>-</sup>
19 <sup>+</sup>	67 <sup>+</sup>		110 <sup>-</sup>
77 <sup>+</sup>	1 <sup>-</sup>	9 <sup>+</sup>	
77 <sup>+</sup>	104 <sup>-</sup>	9 <sup>+</sup>	110 <sup>-</sup>
77 <sup>+</sup>	67 <sup>+</sup>		110 <sup>-</sup>



# After feature hashing

$$\begin{bmatrix} \dots & \mathbf{x}_1 & \dots & \mathbf{x}_2 & \dots & \mathbf{x}_5 & \dots & \mathbf{x}_9 & \dots & \mathbf{x}_{19} & \dots & \mathbf{x}_{67} & \dots & \mathbf{x}_{77} & \dots & \mathbf{x}_{104} & \dots & \mathbf{x}_{110} & \dots \\ -1 & & -1 & & -1 & & & & & & & & & & & & & -1 \\ -1 & & -1 & & & & & 1 & & & & & & & & & & & \\ -1 & & -1 & & & & & 1 & & & & & & & -1 & & & -1 \\ -1 & & & & & & & & & & 1 & & & & & & & -1 \\ & & & & -1 & & & & 1 & & & & & & & & & -1 \\ & & & & & & & 1 & & 1 & & & & & -1 & & & -1 \\ & & & & & & & & 1 & & 1 & & & & & & & -1 \\ -1 & & & & & & & 1 & & & & & & 1 & & & & & \\ & & & & & & & 1 & & & & & & 1 & & -1 & & -1 \\ & & & & & & & & & & 1 & & 1 & & & & & -1 \end{bmatrix}$$

$\text{rating} \sim \text{user} + \text{movie} + \text{has\_cats} + \text{has\_dogs}$

- Average effect for users
- Average effect for movies
- Average effect for movie-related features

$\text{rating} \sim \text{user} + \text{movie} + \text{has\_cats} + \text{has\_dogs}$

- Average effect for users
- Average effect for movies
- Average effect for movie-related features

$\text{rating} \sim \text{user} + \text{movie} + \text{has\_cats} + \text{has\_dogs}$

- Average effect for users
- Average effect for movies
- Average effect for movie-related features

$\text{rating} \sim \text{user} + \text{movie} + \text{has\_cats} + \text{has\_dogs}$

- Average effect for users
- Average effect for movies
- Average effect for movie-related features

# Interactions

$\text{rating} \sim \text{user} + \text{movie} + \text{user} \times (\text{has\_cats} + \text{has\_dogs})$

- Average effect for users
- Average effect for movies
- Average effect for movie-related features
- User correction

# Interactions with feature hashing

$$h(\text{user}=\text{Ashley}) = 2^-$$



$$h(\text{has\_cats}) = 9^+$$



```
1 def interact(h1, h2):  
2     i1, s1 = h1  
3     i2, s2 = h2  
4     return ((i1 ^ i2) * 16777619) & 127, s1 * s2
```



$$h(\text{user}=\text{Ashley} \times \text{has\_cats}) = 81^-$$

## After feature hashing with interactions

user	movie	has_cats	has_dogs	likes_cats	likes_dogs
2 <sup>-</sup>	5 <sup>-</sup>		110 <sup>-</sup>		4 <sup>+</sup>
2 <sup>-</sup>	1 <sup>-</sup>	9 <sup>+</sup>		81 <sup>-</sup>	
2 <sup>-</sup>	104 <sup>-</sup>	9 <sup>+</sup>	110 <sup>-</sup>	81 <sup>-</sup>	4 <sup>+</sup>
2 <sup>-</sup>	67 <sup>+</sup>		110 <sup>-</sup>		4 <sup>+</sup>
19 <sup>+</sup>	5 <sup>-</sup>		110 <sup>-</sup>		71 <sup>-</sup>
19 <sup>+</sup>	104 <sup>-</sup>	9 <sup>+</sup>	110 <sup>-</sup>	110 <sup>+</sup>	71 <sup>-</sup>
19 <sup>+</sup>	67 <sup>+</sup>		110 <sup>-</sup>		71 <sup>-</sup>
77 <sup>+</sup>	1 <sup>-</sup>	9 <sup>+</sup>		12 <sup>+</sup>	
77 <sup>+</sup>	104 <sup>-</sup>	9 <sup>+</sup>	110 <sup>-</sup>	12 <sup>+</sup>	25 <sup>-</sup>
77 <sup>+</sup>	67 <sup>+</sup>		110 <sup>-</sup>		25 <sup>-</sup>



# After feature hashing with interactions

$$\begin{bmatrix} & x_1 & x_2 & x_4 & x_5 & x_9 & x_{12} & x_{19} & x_{25} & x_{67} & x_{71} & x_{77} & x_{81} & x_{104} & x_{110} \\ -1 & & -1 & 1 & -1 & & & & & & & & & & -1 \\ & -1 & -1 & & & 1 & & & & & & & -1 & & \\ & & -1 & 1 & & 1 & & & & & & & -1 & -1 & -1 \\ & & -1 & 1 & & & & & & 1 & & & & & -1 \\ & & & & -1 & & & 1 & & & -1 & & & & -1 \\ & & & & & 1 & & 1 & & & -1 & & & -1 & 0 \\ & & & & & & 1 & & & 1 & -1 & & & & -1 \\ -1 & & & & & 1 & 1 & & & & & 1 & & & \\ & & & & & 1 & 1 & & -1 & & & 1 & & -1 & -1 \\ & & & & & & & & -1 & 1 & & 1 & & & -1 \end{bmatrix}$$

# Pros and cons

## Pros

- Time- and space-efficient → online learning
- Sparsity-preserving
- Implicit handling of missing data

## Cons

- Collisions → statistical performance
- Inverse mapping → interpretability

# Feature hashing

---

## Use case

# Scenario

---

- Business directory service (think Yelp)
- Want to improve relevance of and personalise search results

# Scenario

- Business directory service (think Yelp)
- Want to improve relevance of and personalise search results

## CTR optimisation

Maximise probability of click on top results

## Personalisation

Take into account past user interactions

# Data

user	query	business	click
A	Italian restaurant	Bocca di Lupo	1
A	Italian restaurant	Brasserie Zédel	0
A	Italian restaurant	Emilia's Crafted Pasta	0
A	Italian restaurant	Fucina	1
A	Italian restaurant	Trullo	0

- $3.5 \times 10^8$  observations (250 days)
- $4.5 \times 10^6$  businesses with many (sparse) attributes

$$\text{click} \sim \text{user} + \text{query} + \text{business} + \\ \text{query} \times \text{business} + \text{user} \times \text{business}$$

- Average effects
- Relevance correction
- User correction

$$\text{click} \sim \text{user} + \text{query} + \text{business} + \\ \text{query} \times \text{business} + \text{user} \times \text{business}$$

- Average effects
- Relevance correction
- User correction



$$\text{click} \sim \text{user} + \text{query} + \text{business} + \\ \text{query} \times \text{business} + \text{user} \times \text{business}$$

- Average effects
- Relevance correction
- User correction

$$\text{click} \sim \text{user} + \text{query} + \text{business} + \\ \text{query} \times \text{business} + \text{user} \times \text{business}$$

- Average effects
- Relevance correction
- User correction

- Fixed  $p = 2^{23} \approx 8.4 \times 10^6$
- Regularised logistic regression (Spark ML)
- Grid search over elastic net hyperparameters

# Results

- Estimated  $\approx 5\%$  increase in CTR based on train/test split
- A/B testing showed  $\approx 10\%$  increase in CTR

# Library support

---

# Vowpal Wabbit (VW)

- C++ with Python bindings
- Linear, logistic and Poisson regression with interactions
- Extremely fast and scalable

- Some support for feature hashing:  
    `feature_extraction.FeatureHasher` and  
    `feature_extraction.text.HashingVectorizer`
- No interactions

- Some support for feature hashing in MLlib:  
`ml.feature.FeatureHasher` and `ml.feature.HashingTF`
- PySpark transformers (hashing and interactions)
- VW bindings coming soon to MMLSpark



## Feature hashing is...

- Great for online training on lots of data
- Well-suited for high-cardinality features (sparse  $\mathbf{X}$ )
- Even better with interactions

# Thank you!

If you want to keep in touch...

@ gianluca@campanella.org    ● gcampanella    in gcampanella