

# Feature Hashing for Large Scale Multitask Learning

Weinberger et al. 2009

Seminar aus Data Mining und Maschinellem Lernen



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

U: Votre Apotheke en li  
-10 pilu x 100 mg + Cila  
raison gratuite  
stème de commande sûr



$x$

NEU  
Votre  
Apotheke  
...



NEU  
USER123\_NEU  
Votre  
USER123\_Votre  
Apotheke  
USER123\_Apotheke  
...



$\phi$

1  
0  
-1  
0  
0  
-1  
0  
1  
0  
...

text document (email)

bag of words

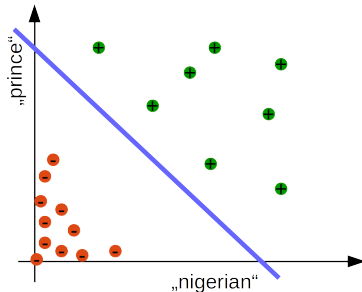
bag of words  
(personalized)

hashed,  
sparse vector  
 $\phi_0(x) + \phi_u(x)$

Task: classify email into spam / not spam

- ▶ using *linear classifier*
  - ▶ classification can be expressed as a dot product of a *feature vector*  $\mathbf{x}$  and a *weight vector*  $\mathbf{w}$

$$c = \langle \mathbf{x}, \mathbf{w} \rangle$$



How to turn a text document into a vector?

naive approach: dictionary

- ▶  $D : \text{words} \rightarrow \mathbb{N}$
- ▶ every index in a feature vector corresponds to a word
  - ▶ one-hot encoding
- ▶ normalized sum over all words in a document

Problem: high dimensionality ( ~40 million)

Size of the feature vector:

- ▶ number of words
- ▶ in every language
- ▶ + every misspelling

⇒ Solution: dimensionality reduction

- ▶ Only consider frequent words?
- ▶ misspellings and unknown brand names are probably good indicators for spam

Reduces the dimensionality by hashing into a space of size  $m$

Advantages:

- ▶ no (original) features are ignored
- ▶ no dictionary
- ▶ outputs fixed-size feature vector
  
- ▶ hash function  $h : \mathbb{N} \rightarrow \{1, \dots, m\}$
- ▶ hash function  $\xi : \mathbb{N} \rightarrow \{\pm 1\}$
- ▶  $\phi_i^{h,\xi}(\mathbf{x}) = \sum_{j:h(j)=i} \xi(j)x_j$



- ▶ hash function  $h : \mathbb{N} \rightarrow \{1, \dots, m\}$
- ▶ hash function  $\xi : \mathbb{N} \rightarrow \{\pm 1\}$
- ▶  $\phi_i^{h,\xi}(\mathbf{x}) = \sum_{j:h(j)=i} \xi(j)x_j$

Pseudocode (assuming  $\mathbf{x}$  is a normalized sum of one-hot encoded vectors):

```
function hashing_vectorizer(features : array of string, N : integer):  
    x := new vector[N]  
    for f in features:  
        h := hash(f)  
        idx := h mod N  
        if  $\xi(f) == 1$ :  
            x[idx] += 1  
        else:  
            x[idx] -= 1  
    return x
```

source: [https://en.wikipedia.org/wiki/Feature\\_hashing](https://en.wikipedia.org/wiki/Feature_hashing)



Properties of this formulation:

- ▶ unbiased

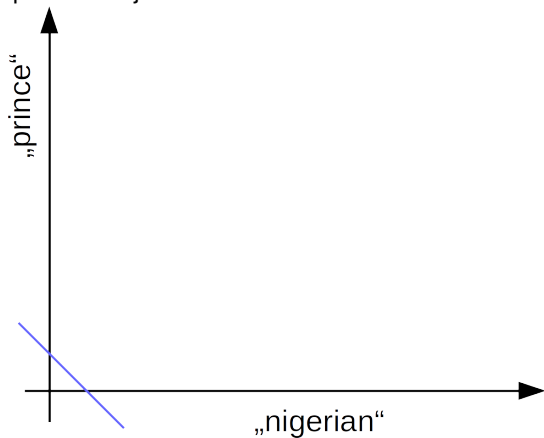
$$\mathbf{E}_\phi[\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle] = \langle \mathbf{x}, \mathbf{x}' \rangle$$

- ▶ variance shrinks with dimensionality  $m$  of the hash-space  
variance rises with big values in  $\mathbf{x}$  or  $\mathbf{x}'$

$$\mathbf{Var}_\phi[\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle] = \frac{1}{m} (\sum_{i \neq j} x_i^2 x_j'^2 + x_i x_i' x_j x_j')$$

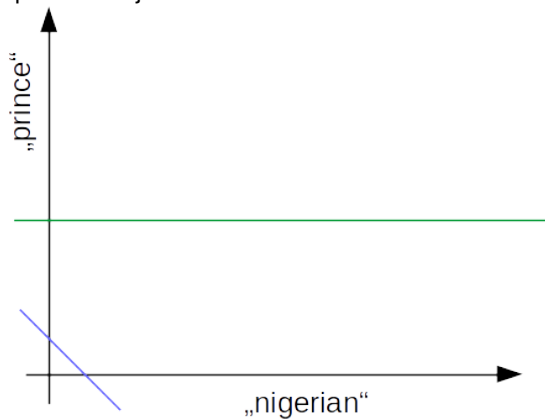
- ▶ length-preserving with high probability
- ▶ the larger a value in  $\mathbf{x}$  or  $\mathbf{x}'$ , the larger the distortion from collisions
- ▶ hashing multiple times
  - ▶ decreases large values in  $\mathbf{x}$  or  $\mathbf{x}'$
  - ▶ decreases sparsity
  - ▶ increases variance

Spam is subjective.

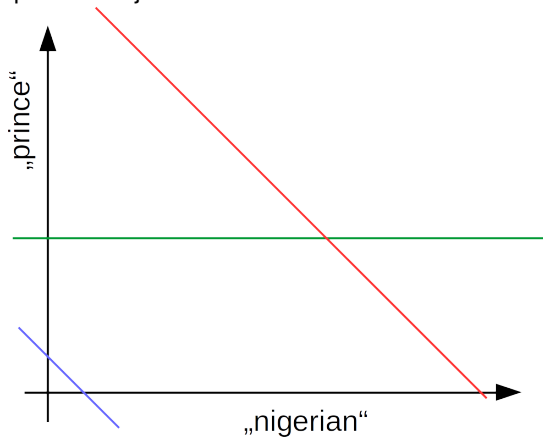




Spam is subjective.



Spam is subjective.



Goal: Learning a number of related tasks (individual spam-filter per user).

Challenges:

- ▶ little training data per user (task)
- ▶ *no* training data for most users
- ▶ learn shared parts of all tasks
- ▶ learn a specialized classifier for each task

⇒ inductive transfer between tasks

Example:

- ▶ **everybody** can agree that “Viagra” indicates spam (shared)
- ▶ **some** people regard “nigerian” + “prince” as spam (task-specific)

- ▶ hash functions  $\phi_0, \dots, \phi_{|U|}$
- ▶ for each user that provided labels:
  - ▶ one local hash function  $\phi_u$
  - ▶ one local predictor  $\mathbf{w}_u$
- ▶ one global hash function  $\phi_0$
- ▶ one local predictor  $\mathbf{w}_0$  shared between all users

All hashed into a common space:

$$\mathbf{w}_h = \phi_0(\mathbf{w}_0) + \sum_{u \in U} \phi_u(\mathbf{w}_u)$$

⇒ assumption: low interference between user hashes

prediction for user  $u$ :  $\langle \phi_0(\mathbf{x}) + \phi_u(\mathbf{x}), \mathbf{w}_h \rangle = \langle \mathbf{x}, \mathbf{w}_0 + \mathbf{w}_u \rangle + \epsilon_i + \epsilon_d$

interference error  $\epsilon_i$

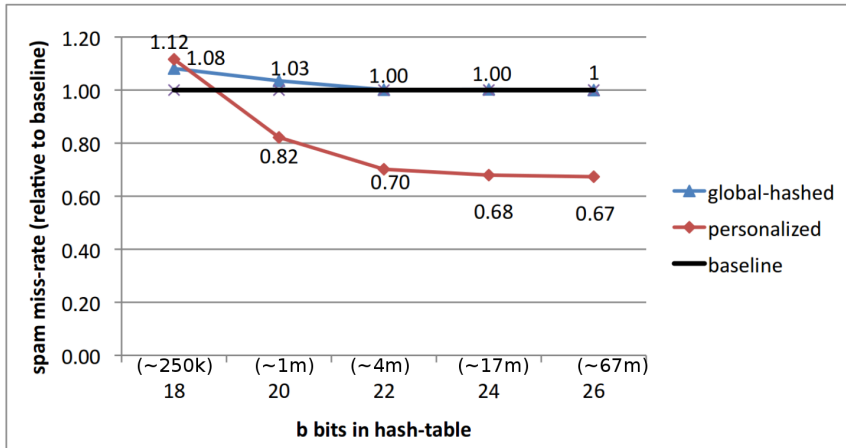
- ▶ caused by collisions between  $\phi_0(\mathbf{x})$  or  $\phi_u(\mathbf{x})$  with hash functions of other users
- ▶ shown to be small with high probability

distortion error  $\epsilon_d$

- ▶ caused by self-collisions of hash functions
- ▶ shown to be small with high probability

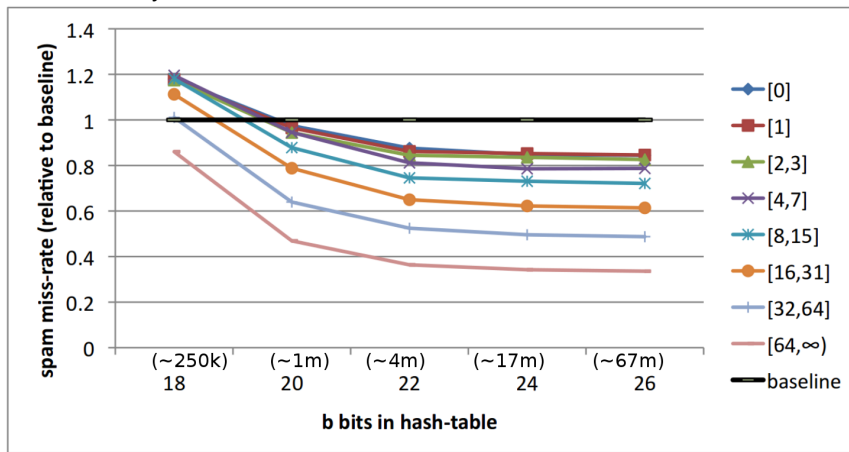
# Results

dimensionality of the raw data set: ~40 million



# Results

dimensionality of the raw data set: ~40 million



# Summary

