# Quora Question Pairs

Can you identify question pairs that have the same intent?

## Naïve Bayes Classifier
**k** 0.64885

**Snowball stemmer** 'having' ↦ 'have'

**Feature engineering**  [question1, question2] ↦ [similar_words, difference]
  Similar words
  Absolute difference of length

**Train a Naïve Bayes Classifier**
  Validation accuracy: 0.65
  Most Informative Feature
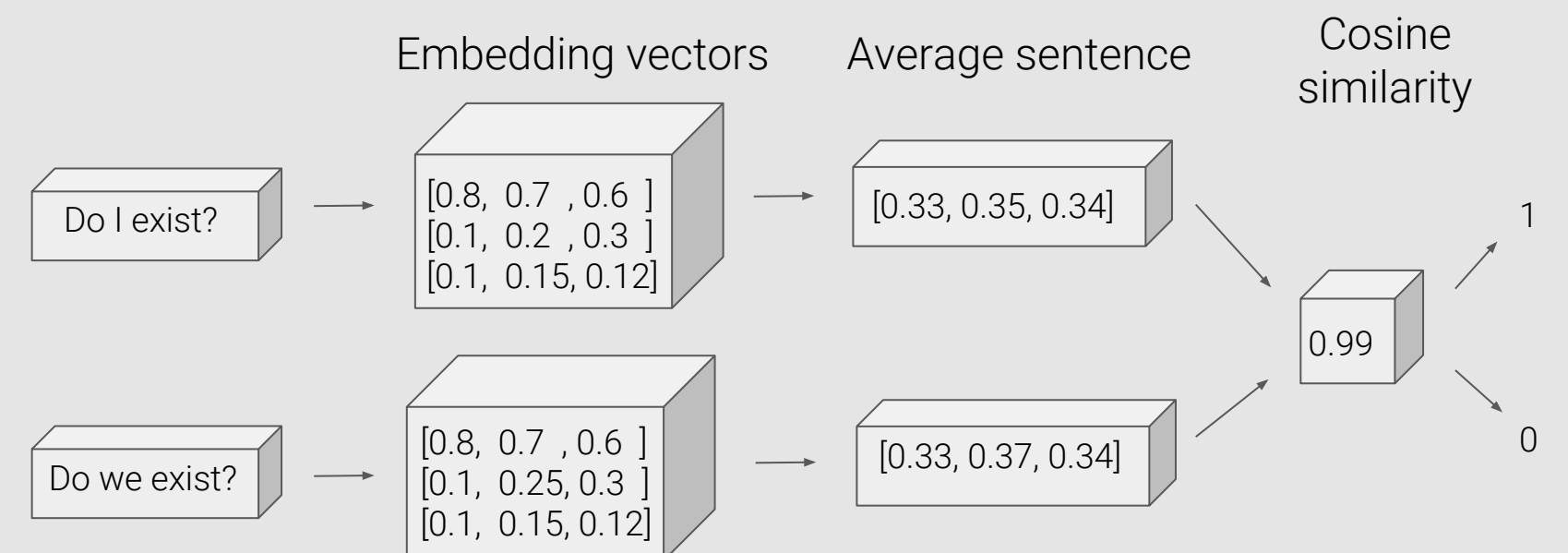    `similar_words = 1  ;  0 : 1 = 1805.6 : 1.0`

## Word2Vec
**k** 0.69281

*[king] - [man] + [woman] ≈ [queen]*

**Hyperparameters**
  Context - window size = 5
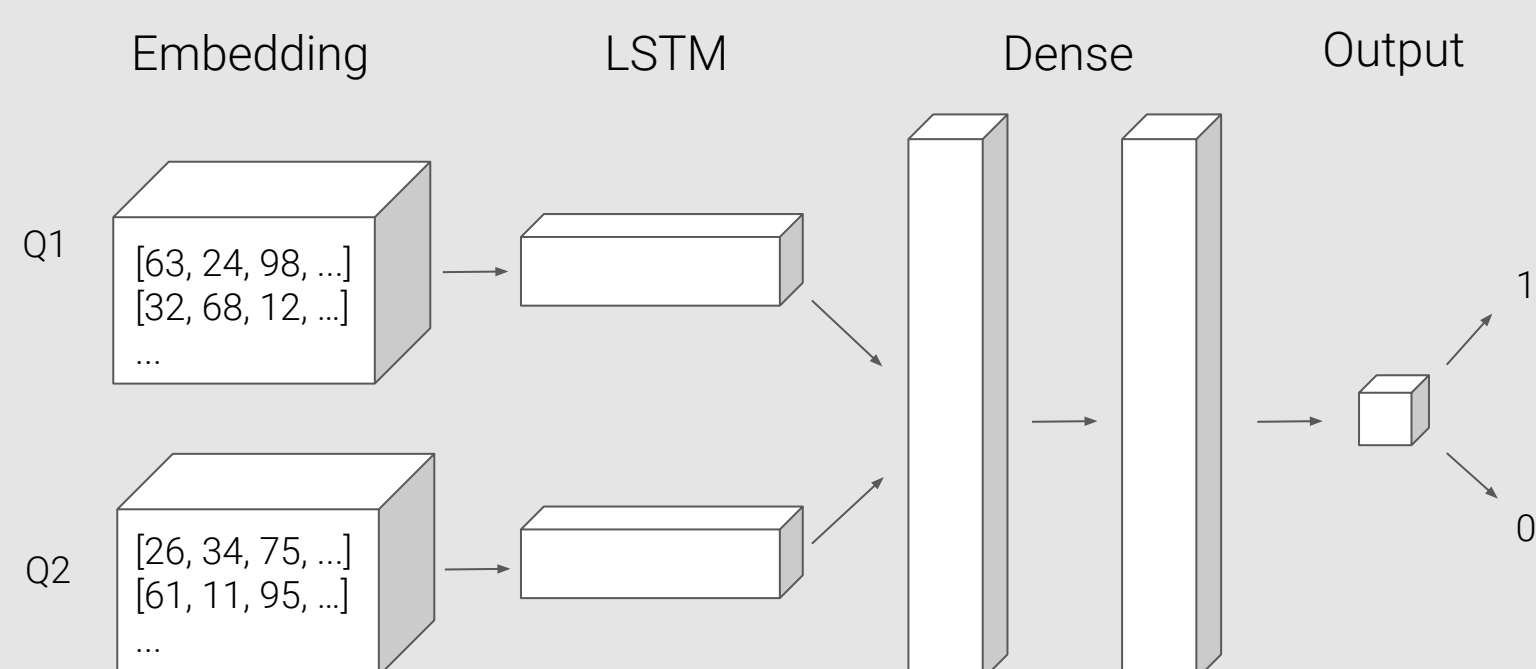**Start from empty model** No transfer learning



Embedding vectors — Average sentence — Cosine similarity

Do I exist? → [0.8, 0.7, 0.6] [0.1, 0.2, 0.3] [0.1, 0.15, 0.12] → [0.33, 0.35, 0.34]

Do we exist? → [0.8, 0.7, 0.6] [0.1, 0.25, 0.3] [0.1, 0.15, 0.12] → [0.33, 0.37, 0.34]

0.99 → 1 / 0
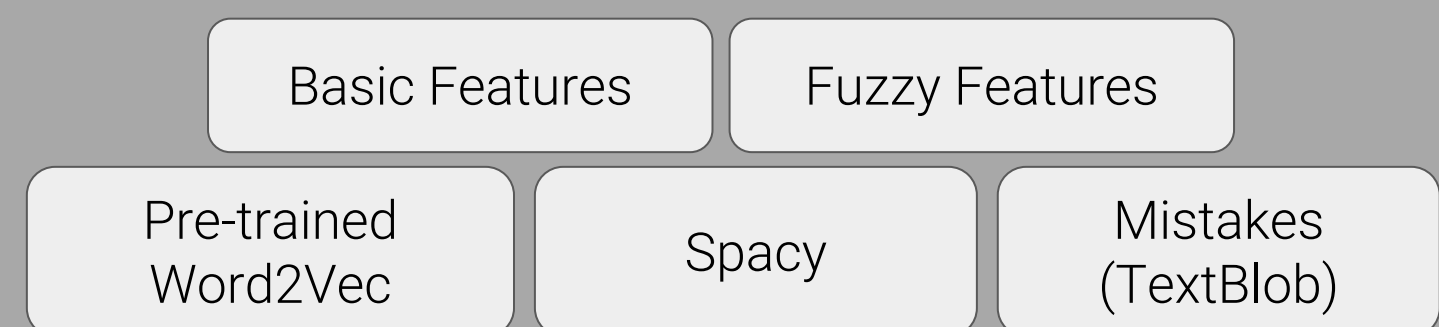
## Deep Learning Model v1.0
**k** 0.78967

| Raw | [What, is, a, Horcrux, ?] |
| Lemma | [what, be, a, horcrux, ?] |
| Tags | [WP, VBZ, DT, NNP, .] |

**+** **Global Vectors (GloVe)**
Source: Common Crawl
840B tokens
2.2M vocab
300d vectors

Embedding — LSTM — Dense — Output

Q1 [63, 24, 98, ...] [32, 68, 12, ...] ...
Q2 [26, 34, 75, ...] [61, 11, 95, ...] ...

→ 1 / 0

## Cleaning and Feature Engineering

Removing the shortest and the longest questions
Adding new features

Basic Features    Fuzzy Features
Pre-trained Word2Vec    Spacy    Mistakes (TextBlob)

## Logistic Regression
**k** 0.81049

L2 penalty
Hyperparameters tuning
  C ∈ [0.001,1000] in a logarithmic step
  Grid search
  Choose C with smallest diff. training-validation accuracy

## dmlc XGBoost
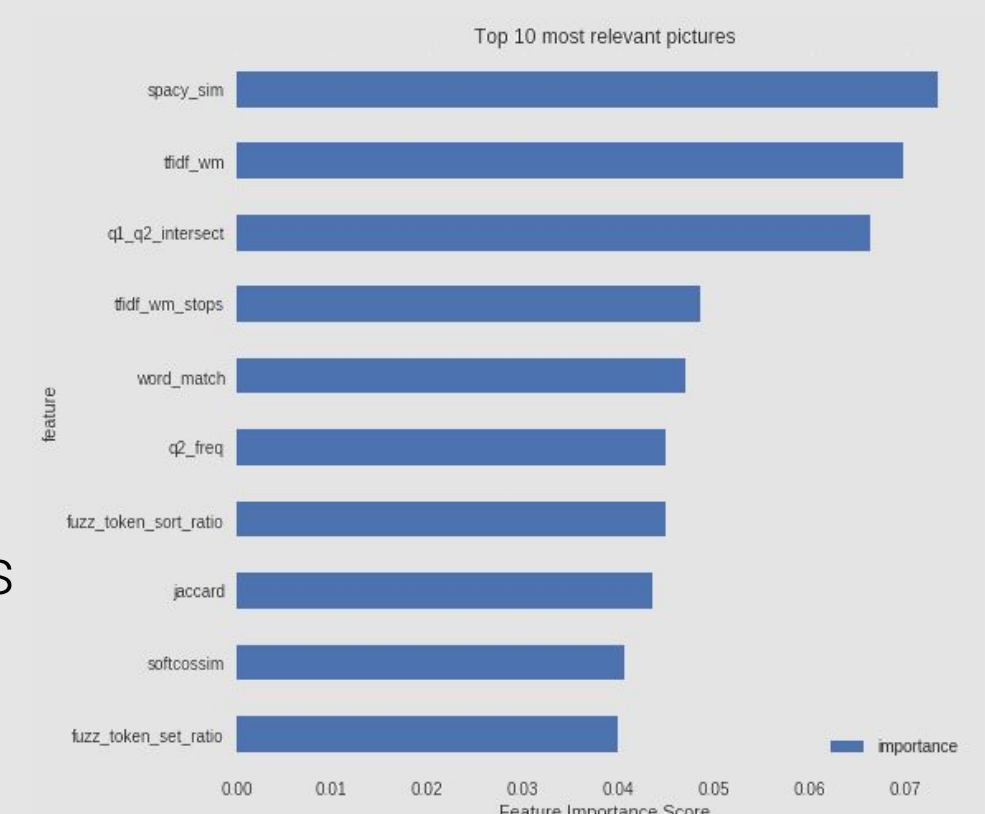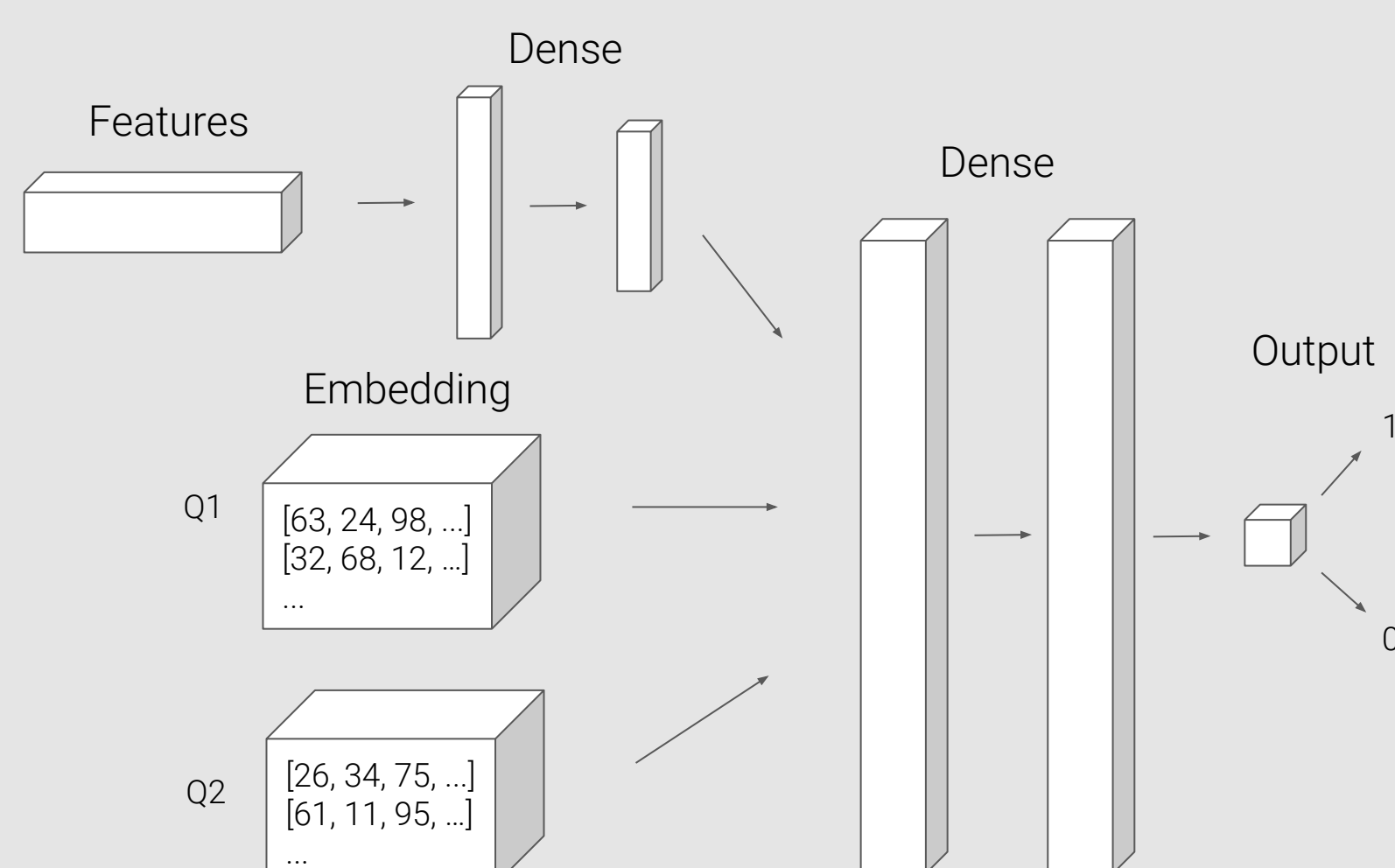**k** 0.81027

**Input**
  Extracted features
  *is_duplicate* label

**Hyperparameters tuning**
  Cross-validation: 10 folds
  Grid search
  Randomized search



Top 10 most relevant pictures

spacy_sim, tfidf_wm, q1_q2_intersect, tfidf_wm_stops, word_match, q2_freq, fuzz_token_sort_ratio, jaccard, softcossim, fuzz_token_set_ratio

Feature Importance Score — importance

## Deep Learning Model v2.0
**k** 0.81929

Features — Dense — Dense — Output

Q1 [63, 24, 98, ...] [32, 68, 12, ...] ...
Q2 [26, 34, 75, ...] [61, 11, 95, ...] ...

Embedding

→ 1 / 0

## Ensemble Learning
**k** 0.85912

**Bagging**
  Less variance
  Fights overfitting

**7 aggregated models (k > 0.79)**
  4x Deep Learning models
  1x Logistic Regression model
  2x XGBoost models

*"The majority cannot be wrong"*
*"Two heads are better than one"*

Each model will have the same voting power:
  [0, 1, 0, 0, 0, 0, 0] → 0
  [1, 1, 1, 1, 1, 0, 0] → 1

kaggle                    Albert **Folch** | Xavier **Moreno**