

pic2kcal: End-to-End Calorie Estimation From Pictures of Food

Robin Ruede, Lukas Frank, Verena Heußner
Karlsruhe Institute of Technology

Abstract

We estimate kcal directly from a picture. It good.

1. Introduction

2. Related Work

There's some other papers like [1]–[3]. Ours is more end to end and also BETTER

3. Dataset Extraction and Preprocessing

3.1. Collection

We collected a dataset from the a popular German recipe website that contains ingredient lists, cooking instructions, and pictures of the resulting meals. The recipes are from many different cuisines and also include things like cakes, cocktails, and other. Most recipes have at least one picture. Pictures can come both from the original author of the recipe, as well as from third parties. Some of the pictures are of a single plate of food, others are for example of a whole casserole. We do not have any information about whether a picture contains a single portion. Around 10% of recipes contain a user-given value for how many calories per portion the recipe has.

3.2. Matching / Preprocessing

Since the dataset only has user-given calorie information for a small part of the data and doesn't include any details regarding the macronutrient composition, and since the user given information is often inaccurate [todo: compare to ours?], we decided to match the list of ingredients against a database of nutritional values to sum up the proportions of macronutrients as well as the total calories.

To facilitate this, we collected a secondary dataset from a German website of nutritional values. The website contains values for the amount of fat, protein, and carbohydrates in grams per 100g of product. Additionally, it contains user-sourced food amounts like "1 medium-sized apple = 130g". The data is partially sourced from the USDA Food Composition Database [4], and partially crowd sourced from

manufacturer-given data of specific products.

Matching the recipe ingredients to the nutritional database has two main problems.

Firstly, the given ingredient name often includes information that is not relevant to the product itself, but rather to its preparation or visual qualities. These additional text snippets are hard to separate from information that is relevant. For example: 3 onions, diced and 3 onions, in slices refer to the same product, while 500g pasta, cooked and 500g pasta, raw vary significantly in their caloric density. We experimented with three approaches [reformulate] to solve this problem. Firstly, we tried simply matching the ingredient to the nearest ingredient based on character edit distance (Levenshtein distance). This resulted in very bad matchings because of missing handling of synonyms and the above issue. To solve this we tokenize the ingredient name to words, embed each word to a vector with Word2Vec or FastText [5], and then use average the word vectors to get a ingredient vector. This is the same method as used in the fasttext library for extracting sentence vectors. This still lead to unsatisfactory results, since each word in the ingredient name has the same weight, even though some specify less important details. For example in "red onion" vs "red apple", the word "red" is much less important than "onion" and "apple". We got the best result by using the Google Universal Sentence Encoder [6], which creates 512-dimensional embeddings of any amount of text. We find the best matches for a ingredient by comparing the embedding of user-given free text from the recipe to the embeddings for all food items for which we have nutritional data using the cosine distance, and then try to find a conversion for the given amount to a normalized gram or milliliter amount.

The second problem is matching the amounts. For ingredients given in grams this is trivial, but for many items the recipe authors use other units of measure like e.g. can, piece, tablespoon, "some", "2 large X", "salt 'by taste' ". Since spices usually have little impact on the nutritional values, we exclude ingredients that are "by taste" and similar. For the other amounts, we match the unit name (like tablespoon or "medium large") exactly and multiply it with the given amount. We also add some special cases like match-

ing “can” to “can (drained weight)” and similar.

The amount matching is applied to all possible ingredient matches that are similar by more than 84% (measured by cosine distance) [why lol] in decending order, or to the single closest ingredient if there is no match more accurate than 84%.

If the amount matching fails, the ingredient is marked as unmatched. If a recipe has at least one unmatched ingredient, it is discarded. With this method we have full nutritional information for around 38% of all recipes. This could be improved with further tweaking.

3.3. Dataset Statistics

- statistics about dataset
- number of pictures per recipe
- number of recipes
- most common ingredients
- cake bias

4. Models

- (SmoothL1Loss + weight * BCE)
- What we predict (kcal, fat, etc, ings)

5. Experiments

- todo: actually use test data set

| Method | kcal relative error |
|-----------------------|---------------------|
| baseline | 0.464 |
| ours (kcal only) | 0.361 |
| ours (w/ macros) | 0.352 |
| ours (w/ macros+ings) | 0.328 |

Table 1. Results per 100g. Note that multitask learning improves performance.

6. Results

Our results can be seen in tbl. 1. Example outputs can be seen in fig. 1.

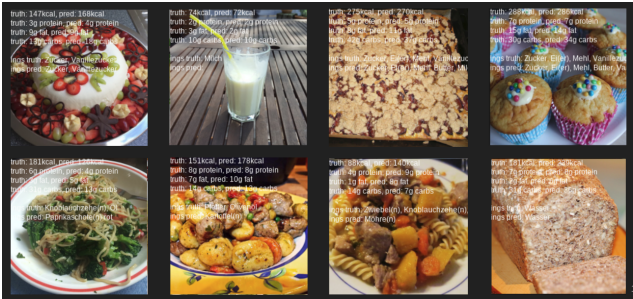


Figure 1. Some example results, showing predicted calories, fat, protein, carbohydrates and ingredients.

References

[1] M. Chokr and S. Elbassuoni, “Calories Prediction from Food Images,” in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, 2017, pp. 4664–4669.

[2] T. Ege and K. Yanai, “Image-Based Food Calorie Estimation Using Knowledge on Food Categories, Ingredients and Cooking Directions,” in *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*, 2017, pp. 367–375.

[3] A. Romero, X. Giro-i-Nieto, M. Drozdal, and A. Salvador, “Inverse Cooking: Recipe Generation from Food Images,” Dec. 2018.

[4] *USDA Food Composition Databases*.

[5] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching Word Vectors with Subword Information,” *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 135–146, Dec. 2017.

[6] D. Cer *et al.*, “Universal Sentence Encoder,” in *In submission to: EMNLP demonstration*, 2018.