

# Beyond Metadata for BBC iPlayer: an autoencoder-driven approach for embeddings generation in content similarity recommendation

Simone Spaccarotella

September 2024

## Contents

<b>1</b>	<b>Introduction and background</b>	<b>3</b>
<b>2</b>	<b>Outline of the issue or opportunity and the business problem to be solved</b>	<b>4</b>
<b>3</b>	<b>Methods and justification</b>	<b>6</b>
3.1	Data pre-processing . . . . .	6
3.2	Modelling and regularisation . . . . .	6
3.3	Content similarity . . . . .	7
3.4	Tools and frameworks . . . . .	7
<b>4</b>	<b>Scope of the project and Key Performance Indicators</b>	<b>9</b>
<b>5</b>	<b>Data selection, collection and pre-processing</b>	<b>10</b>
<b>6</b>	<b>Survey of potential alternatives</b>	<b>12</b>
<b>7</b>	<b>Implementation and performance metrics</b>	<b>14</b>
<b>8</b>	<b>Discussion and conclusions</b>	<b>17</b>
8.1	Results . . . . .	17
8.2	Summary of findings and recommendations . . . . .	17
8.3	Implications . . . . .	18
8.4	Caveats and limitations . . . . .	18



# 1 Introduction and background

I am a Software Engineer at the BBC, Team Lead for the Sounds web team, and I have been training as a Data Scientist, working in attachment to the iPlayer Recommendation team.

I built a machine learning model pipeline that generates content-to-content (C2C) similarity recommendations of video-on-demand (VOD) for the “More Like This” section on BBC iPlayer [3]. This project is relevant to me because I have been crossing paths with the world of recommendations multiple times during my career at the BBC, which sparked my interest. I had a tangent encounter in 2015 while working for a team that built an initial recommender for BBC News and an API to provide recommendations using 3rd party engines. I also produced and presented a talk for a Hack Day. The talk was called “Recommendation Assumptions” [17], and it was about types of recommendations and contextual external factors affecting them. Until now, when I was able to finally put my knowledge into practice with an actual project on real data.

The BBC is a well-known British broadcaster that constantly evolves to remain relevant to its audience. Its mission is to inform, educate, and entertain, and it operates within the boundaries set by the Royal Charter [7]. The current media landscape requires the BBC to deliver digital-first content relevant to the audience. This transformation involves investments in data and personalised services, not to mention a certain revolution in generative machine learning modelling that is keeping everyone busy.

## 2 Outline of the issue or opportunity and the business problem to be solved

The BBC produces and stores vast amounts of data for its content, surfaced by countless services and APIs. One of the top priorities for the BBC is to increase the usage across the business of **Passport** [8], an internal BBC system that provides a richer set of metadata annotations for multimodal content (audio, video and text). The usage in production is low, and its BBC-wide adoption would make access to metadata consistent, remove duplications, and reduce effort and costs.

Furthermore, the similarity score of the current C2C recommender is directly proportional to the number of values in common between any pairs of items on a per-feature basis. However, the commonality is calculated with exact string equality, ignoring any relationship between different categorical values expressing a similar concept (e.g. “comedy” and “stand-up comedy”).

Moreover, the limited number and type of tags cannot adequately describe the content. At the same time, the skewedness of the data distribution and a lack of pre-processing, shifts the recommendations towards the most popular categories.

Lastly, each similarity score is multiplied by a hardcoded weight that modulates the importance of a feature, but it doesn’t solve the polarising effect of a skewed distribution. Unfortunately, because these are hyperparameters and not learned weights, the model can’t improve its performances by minimising them against a cost function.

To address these issues, the aim of this project was:

- **To improve the quality of the C2C similarity recommendations.** By using in input a richer set of metadata that better describes the content, and by reducing the high-dimensional data to a lower-dimensional latent manifold, the model would be able to generate embeddings that could improve the quality of the recommendations, by mapping the item similarity problem to a geometric distance calculation between vectors in a multidimensional Euclidean space.
- **To build a general solution that can be applied to multimodal content, reducing the costs.** I built a C2C recommender for iPlayer, using the Passport dataset to make the solution general so that it could be applied to any BBC content, which shares the same set of standard tags.

- **To build a foundational item-embeddings generator.** Content-based recommenders use item metadata. This project provided an immediate solution for non-personalised C2C recommendations that solely rely on them. It also provided a foundational basis for personalised recommenders that could benefit from using content metadata embeddings combined with other data like user interactions.

## 3 Methods and justification

### 3.1 Data pre-processing

I used **one-hot encoding** to transform the categorical features (i.e. the metadata annotations) into a numerical vector. It is a simple yet effective encoding method and is perfect for transforming nominal categoricals because it doesn't introduce any ranking or arithmetic relationship among the encoded values. The downside of this approach is that it generates high-dimensional sparse arrays, introducing the so-called *curse of dimensionality* problem. Nonetheless, this was an accepted drawback that was managed in the modelling phase.

### 3.2 Modelling and regularisation

I trained an **autoencoder** [1, 15] to learn the Passport tags' latent features and reduce the encoded vectors' size. The autoencoder is an encoder-decoder neural network, a self-supervised model capable of capturing non-linearity from the data. I used the "undercomplete" variant, which constrains the number of nodes in the hidden layers, creating a "bottleneck" of information flow through the network. This bottleneck is a form of *regularisation* that forces the model to learn latent attributes from the input while reconstructing it with minimal loss. Ultimately, it helps prevent the model from overfitting the training dataset by indexing it like a caching layer.

I extracted the trained *encoder* segment of the network to compress the one-hot encoded high-dimensional sparse array into a lower-dimensional and denser representation called *embedding* [10]. This technique solved the curse of dimensionality and the data sparsity problems and improved the calculation complexity and the quality of the recommendations at inference time.

To improve and assess the ability of the model to *generalise* on unseen data, I randomly shuffled the dataset and split it into three chunks: training, validation and test. I used **hyperparameter tuning** to find the best model parameters that minimised the cost function and used the validation set to regularise the model with **early stopping**. This technique monitored the reconstruction loss on an out-of-sample dataset, allowing the model to stop training within a set "patience" threshold after reaching a local minimum on the validation error. The test set was finally used to assess the model's performance using the best set of parameters.

I used **dropout** to further regularise the model and make it robust to small changes in the input, and **data augmentation**, by including in the training data the episodes that shared the same tags with their parent pro-

gramme. **Weight decay** and **batch normalisation** were tested during hyperparameter tuning and discarded for poor performance.

### 3.3 Content similarity

Content similarity was calculated with the **cosine** of the angle  $\theta$  between each pair of embeddings [11]. This metric is insensitive to the magnitude of the vectors. Because high-frequency values tend to have a larger magnitude, it mitigates the impact of popularity in the similarity calculation. One-hot encoded vectors lack meaningful relations between them. They represent unit vectors bound in the “positive quadrant” of a Cartesian coordinate system for a multidimensional Euclidean space. Because each pair can only have a finite number of angles, the cosine similarity will also assume a finite number of discrete values between 0 and 1, causing information loss. Ideally, we would expect the similarity score to assume a continuous value bound between -1 and 1, and this is only possible if the angle  $\theta$  of any vector pair can assume a value between 0 and 360 (i.e.  $0\pi$  and  $2\pi$ ), hence the use of embeddings.

### 3.4 Tools and frameworks

The entire project was written in **Python**. It is the *de facto* programming language for data science and machine learning tasks. Python has an established, diverse and well-documented ecosystem of external libraries and frameworks that facilitated the job, and it is also the language of choice at the BBC.

I used **Pandas** only for tabular data manipulation to generate and store the one-hot encoded vectors. Unfortunately, using it for exploratory data analysis (EDA) wasn’t possible because the iPlayer catalogue had roughly one year’s worth of data, which didn’t fit in memory, causing Pandas to crash. Therefore, I used **Dask**, a library capable of running out-of-memory and parallel execution for faster processing on single-node machines and distributed computing on multi-node machines while using the familiar Pandas API.

I used **TensorFlow** and **Keras** for modelling to build and train the encoder-decoder neural network architecture. **Keras Tuner** for hyperparameters tuning. I also used **Scikit-learn** but not for modelling. It provided utility functions for the dataset splitting and the cosine similarity calculation, and I was already familiar with its API.

I used a combination of **Matplotlib** and **Seaborn** for visualisation and **rdflib** to fetch and parse the RDF documents from the BBC Ontology to

extract the labels needed to visualise the recommendations for testing purposes. Worth also mentioning the use of **pytest** for unit testing and **black** for PEP 8 code compliance and formatting. Finally, I used **Jupyter Lab** to edit the project, **git** for code versioning, **GitHub** as a remote code repository and for collaboration, and **AWS Sagemaker** to run the pipeline on more capable virtual machines, especially during hyperparameter tuning.



## 4 Scope of the project and Key Performance Indicators

The scope for this project was to build an end-to-end machine learning solution, able to produce non-personalised content-to-content similarity recommendations, using Passport metadata tags as input.

The minimum viable outcome was to produce recommendations comparable with the ones currently in production, while a desired outcome was to increase user engagement. The integration with Passport would make this solution general, and applicable to multimodal BBC content. If adopted by  $N$  BBC products with a total cost  $C$ , it could generate considerable savings, with an approximate reduction in costs by a factor of  $\frac{C}{N}$ .

Comparability was a qualitative and subjective key performance indicator (KPI) that served as a compass, indicating whether the project was progressing towards the right direction. It was evaluated by several technical and non-technical stakeholders, with a diverse domain knowledge and background. I built a rudimentary visualisation tool that rendered the title, image, description and metadata of the seed programme and the top-K C2C recommendations. The people involved, gave their subjective feedback on their perceived level of similarity of the output, testing edge cases, common use cases with an expected outcome and sensitive recommendations like content recommendations for children accounts. They also discussed anomalies and/or surprising results.

To measure user engagement though, the solution needed to be A/B tested in production, with real data. Unfortunately, there were too many moving parts outside my control that needed to happen, for me to build a production-worthy version for A/B testing. Adopting this as a KPI would have delayed the project, increasing the odds of failure. To mitigate this risk, I had to decouple it from the success of the project.

I defined a hypothesis testable offline, within my area of control. The hypothesis stated that long-term user engagement is not just about accuracy. It can be affected by increasing diversity in the recommended content. A diverse set of recommendations generates new and unexpected results, with an increase in surprise and serendipity, pushing the user away from boredom. This was an untested hypothesis, but grounded in active research on the topic, such as [13] and [9], that supported the initial statement, and made it less far-fetched. For this reason, I defined a proxy metric that could measure diversity offline, pending a future A/B testing for the validation of the hypothesis, which was pushed out of scope.

## 5 Data selection, collection and pre-processing

News and Sports articles, iPlayer videos on-demand, Sounds podcasts etc. are annotated with the so-called Passport tags. These tags describe any content produced by the BBC and can be used for retrieval (search) and filtering (recommendations). They can be applied either manually by an editorial team with domain knowledge, or semi-automatically by machine learning algorithms with human supervision.

Passport tags are distributed across the BBC via the universal content exposure and delivery (UCED) system, a self-service metadata delivery platform that exposes data as a document stream for products to integrate with. This platform provides different types of *consumers* such as REST API, AWS S3 bucket, etc. Passport documents are JSON objects that contain a property called “**taggings**”, an array of objects representing the metadata annotations. These objects are described by two properties: “**predicate**” and “**value**”. They represent the name and the value of a tag, and are expressed as URL-formatted strings, except for dates. The predicate is a class of the BBC Ontology [2] while the value can be a date or an entity defined as an RDF [20, 18] document, accessible in Turtle format [19] via the BBC Things API [4, 5, 6]. These entities are linked to each other and/or to external resources, and are described by attributes and relationships, giving the data a graph structure.

I decided not to integrate with UCED during development, but to use batches of Passport files, manually collected and stored on a local folder. This was a trade-off that allowed me to train the model with real data, while still keeping the costs down. In addition, because the resources needed to be created on two AWS accounts, I didn’t want to pass the burden of maintenance to the team that owned them, without having tested the feasibility of the solution first.

Content metadata does not constitute personal data and therefore is not subject to the UK GDPR [12]. Nonetheless, this data is encrypted at rest and in transit by default. For this reason, no further actions were required during storage and processing.

I chose to use Passport because it provides a set of tags shared across all content produced by the BBC, making this a general solution that reduces duplications and ultimately costs. Passport provides a flexible and rich set of tags to describe any type of content. Annotations can describe canonical information such as *genre* and *format*, but also things like the *contributor* featuring in the programme, the *narrative theme*, the *editorial tone*, what the content is *about* or what relevant entities are *mentioned*, etc.

During pre-processing, a list of JSON files was loaded into the pipeline

and the tags extrapolated into a dictionary data structure. The key of the dictionary was the programme ID (called *PID*) and the value was another dictionary describing the annotations. A programme can be tagged with the same predicate multiple times, if it has different values, while the same value (e.g. “Music”) can be used by multiple predicates (e.g. **about** or **genre**).

The dictionary was transformed in a Pandas *Dataframe*, where the rows represented the programmes and the columns the tags. I used a MultiIndex [16] for the columns because I needed to keep the duplicate values (2nd-level index) across the predicates (1st-level index). I then populated the cells of the *Dataframe* with the value “1” if the programme was annotated with the corresponding tag, or “0” otherwise, generating one-hot encoded arrays. I initially started with a vectorisation approach known to be performant, using the “**get\_dummies**” Pandas function. Surprisingly, it was slower and less scalable of the solution that I adopted in the end.

A source of bias in the dataset was the “**mentions**” tag. This type of annotation is automatically generated by an algorithm that extracts terms deemed important, appearing in the text of an article or the transcript of an audio/video content. If something is “mentioned”, doesn’t necessarily describe what the content is about, because of the intrinsic ambiguities of natural languages. Figure of speech devices such as metaphors, analogies, allegories, etc., alter the meaning of a sentence for stylistic effect, and can misrepresent the main topic. For example, if the phrase “being over the moon” is mentioned by someone delighted about something unrelated to the topic of “space” and “universe”, extracting the term “moon” as a descriptor, could mislead the representation. To mitigate this source of bias, I dropped the tag in favour of “**about**”, another tag that describes what the content is really about. This is manually added by a team of editorials who are trained to tag the content with relevant and topical annotations.

Generating embeddings of one-hot encoded vectors with the same size used in training but with unseen tags, leads to unpredictable errors. This is because the encoding is positional, and the combination of 1 and 0 learned by the model belongs to the tags seen during training. So, I decided to drop the new tags and encode only the ones the model was trained on, while padding the rest with zeros, pending retraining to capture the new information. Also, some programmes didn’t have any annotations entirely. Adding them to the training data, created a group of entries with all zeros. If enough observations shared this uninformative characteristic, it could have been picked up by the model. I decided to drop these programmes and include only the ones with at least one annotation.

## 6 Survey of potential alternatives

I initially considered to use **clustering** on the one-hot encoded features, to group similar items together. The model would have returned the items belonging to the same cluster of the one considered for similarity. But, I didn't know how many clusters there could be in the data, and I didn't need to. I could have used a density-based technique to autodiscover them, but even in that case, some of the clusters could have had less than K items in a top-K similarity scenario. I could have returned the items belonging to the nearest cluster if needed, but because of this uncertainty, clustering wasn't very useful in this use case.

The geometric interpretation of similarity between two items, is the distance in space between the two vectors representing them. Any item has a degree of similarity with all the others, and clustering was just a coarse-grained discretisation of that concept. I needed a more granular approach, where every item could be compared with anyone else. In geometric terms, I had to calculate the **pairwise distance** between all vectors, given a metric. So, clustering was discarded as a candidate option.

One-hot encoding doesn't use any spatial proximity information to transform the categorical features into ones and zeros. It's a transformation process that pivots the unique values of each original feature, to be the new variables of the transformed vector. If we project these raw vectors in a multidimensional space, we wouldn't be able to use their relative position to each other as a measure of similarity. Moreover, the high dimensionality of the vectors would have increased the computational complexity

To calculate the pairwise distance efficiently, and to produce a meaningful representation of similarity, the vectors needed to be in a denser and lower dimensional space, a manifold embedded into the original high-dimensional ambient space.

Dimensionality reduction techniques such as **principal component analysis (PCA)**, **independent component analysis (ICA)** or **linear discriminant analysis (LDA)** were considered, but there was a problem with them too. Their job is to find a linear projection of the data, but this is a strong assumption that misses important non-linear structures. I discarded the idea of using PCA (which I was more familiar with) but not the idea entirely of using dimensionality reduction. I just needed a non-linear approach, and I turned my attention to manifold learning.

Before introducing the chosen approach in the next section, and discussing the pros and cons, I'd like to describe the alternative pre-processing step I also considered, to generate vectors that didn't suffer from the curse of dimensionality to begin with.

**Hashing** is a non-invertible transformation that can be used for feature reduction. It can generate smaller vectors than one-hot encoding, but I didn't adopt it because hashing has more hyperparameters, which increases its complexity. Most importantly, it introduces the *collision* problem, where two distinct inputs can be mapped to the same index in the same target domain. This issue could be mitigated by choosing the latest and strongest algorithm to reduce the likelihood of collisions, but the trade-off was too computationally expensive for pre-processing.

## 7 Implementation and performance metrics

I used an undercomplete autoencoder to compress the high dimensional one-hot encoded vectors to a lower dimensional embedding representation. This technique captured non-linearity by learning the underlying latent structure of the data, which was key to represent the original Passport tags in a geometrical space, exploiting local proximity as a measure of similarity.

Because the autoencoder has a symmetrical architecture between the *encoder* and the *decoder*, the input and output layers had the same dimensions: 8982 nodes. This number corresponded to the size of the one-hot encoded array.

Some of the hyperparameters were decided based on the nature of the problem. The input data was a tensor of zeros and ones, and the main objective of the network was to reconstruct the output with minimal loss. I used the *Sigmoid* activation function for the output layer, and *binary cross-entropy* as loss function, to allow the network to push the values of the reconstructed output as close as possible to 0 and 1. I didn't use the *SoftMax* activation function for the output layer, because each single node needed to be able to assume those values. I didn't use any residual-based loss function either - such as MSE or MAE - because they don't have a very steep slope for 0/1 values, to allow the gradient to move fast enough. Which is why the negative log loss was the best choice, because it massively penalises huge differences between  $y$  and  $\hat{y}$ , with a natural logarithmic progression. I used the *Rectified Linear Unit (ReLU)* as activation function for the hidden layers, while testing other variants like *Leaky ReLU (LReLU)* and *Parametric ReLU (PReLU)*, but without any relevant improvements in performance. Some other hyperparameters were:

- the optimizer: Adam
- number of epochs: 100
- batch size: 300
- dropout rate: 0.2
- early stopping patience: 10
- early stopping monitoring: binary cross entropy on validation set
- data split: 80% training, 10% validation, 10% test

The remaining hyperparameters were selected through a final round of optimisations. I decided to use a "Bandit-based" approach called *Hyperband* [14], which improves upon *Random Search* by running fewer epochs on the randomly sampled set of parameters, and move on to the next stage by only testing the best performing ones, returning a ranked list of the best hyperparameter sets. The hyperparameters considered by Hyperband were:

- the number of hidden layers
- the embedding size
- the learning rate
- whether to use dropout
- whether to use batch normalisation

The number of hidden layers was a positive odd integer  $N > 0$ , where  $N - 1$  layers were used for the encoder and the the decoder, and one for the bottleneck in the middle. If the value was 1, the network would only have the bottleneck.

The number of nodes per layer was a positive integer  $N > 0$ . Unfortunately, the number of combinations was too high to be meaningfully optimised. To reduce the complexity, I decided to couple the number of nodes per layer as a progression of integer divisions by 2. Each hidden layer in the encoder had half the number of nodes compared to the previous one, and double the amount of the successive one (and *vice versa* in the decoder), except for the layers close to the input and output. In that case, they could assume any of the set values, depending on the number of hidden layers and embeddings size. For example, because I had 8982 nodes in input, the progression of layer dimensions was: 4491, 2245, 1122, 561, 280, 140 and so on. The embedding size was also bound to assume one of these values, and would also impact the maximum number of hidden layers allowed. Therefore, a network with 5 hidden layers and an embedding size of 280 would have had the following configuration: 8982 -> 1122 -> 561 -> 280 -> 561 -> 1122 -> 8982. While a network with 3 hidden layers but an embedding size of 2245 would have had the following configuration: 8982 -> 4491 -> 2245 -> 4491 -> 8982, representing the maximum extension of the progression.

The strength of this approach was the efficient use of space, and the fast inference time. The space scaled linearly with respect to the size of the input to store the embeddings, and quadratically to store the similarity scores. Hyperparameter tuning took quite some time, but this drawback

was compensated by the fact that training was performed offline, while the recommendations were cached and served instantly. In addition, re-training could be scheduled to cover for the new programmes added to the iPlayer catalogue, and the unavailable ones being removed. The main weakness was the model’s interpretability. The autoencoder is a *black box* by definition, and the use of embeddings to calculate the similarity just made the problem worse. It was practically impossible to interpret which of the tags influenced the ranking in the top-K list, by untangling the weights and biases of the neural network. It was also difficult to provide explanations using model-agnostic techniques like SHAP, because the recommendations were calculated by some distance metric, applied on humanly meaningless embeddings, that needed to be linked to the original input, that in turn needed to be decoded back to the original tags.



## 8 Discussion and conclusions

### 8.1 Results

This is a general solution that works with any content that uses Passport tags, and could serve recommendations for multiple BBC products. The adoption of this solution will reduce effort, duplication of code and data, and as a consequence, costs. I shared the findings with the stakeholders, explaining the main benefits and showing the results using the visualisation tool I built. I presented it once to the data scientists and engineers of the iPlayer recommendations team I worked with, and another time to the team in charge of the non-personalised recommendations for the entire BBC. The feedback was positive in both cases and we discussed at length how to move forward with this project, including what to build as an MVP to finally A/B test the solution.

### 8.2 Summary of findings and recommendations

The results were perfectly aligned with the initial objectives and measure of success set at the beginning of the project. My recommendation was to build an initial minimum viable product (MVP), consisting of a Sagemaker pipeline built on the AWS development account, that ingested batch Passport tags. This would allow us to break down the engineering effort, and spot any blockers/challenges that need to be addressed as early as possible, so that we can correct them and/or reconsider some of the assumptions ahead of the production build.

We would need to build two pipelines, one for training and one for inference to generate the embeddings and the similarity scores. The embeddings and the similarities score need to be cached to improve performances. The second stage of this approach would require the integration with UCED to fetch real-time data automatically.

If this solution is viable and passes the A/B test, it could also be used to generate embeddings for other personalised recommenders that use item metadata in conjunction with user interactions and/or contextual data such as day and time of interaction, location, device used, etc.

The project could be further expanded by exploiting the graph nature of the data using graph neural network (GNN), and in particular, a graph autoencoder (GAE) to learn meaningful representation of the graph data, capturing the topological structure and the node content. This could improve upon the current autoencoder, that flattens the graph structure in a list of tags and relies on the positional encoding of these tags to generate the

embeddings. This effort will require further research and prototyping.

### 8.3 Implications

The project presented a unique opportunity for me to work on an end-to-end machine learning pipeline from data preprocessing to inference, practicing my technical skills, building a real neural network, learning about embedding techniques, and content-based recommendation systems. The positive feedback from stakeholders has reinforced my professional confidence and provided invaluable experience in presenting data-driven solutions to a business audience.

For my colleagues and the team, this project has established a replicable framework for C2C similarity recommendations that can be adapted to other BBC products. The modularity of the solution enables flexibility in extending it to multimodal content, enhancing the potential for collaborative developments across departments. This can promote knowledge sharing and foster a data-centric approach to problem-solving within the wider team, as members can leverage this solution to address similar business problems.

For stakeholders and the business, the project presents a scalable solution to reduce data redundancies, decrease maintenance overhead, and potentially reduce costs associated with content recommendation systems. The approach taken provides a standardized method for generating “More Like This” suggestions, thereby improving user engagement on non-personalised content, with a consistent experience across the BBC portfolio. Moreover, the project’s adaptability encourages strategic, data-driven content management across the organization, supporting future initiatives with robust foundations for content similarity and recommendation. Overall, this project not only aligns with business goals of optimizing resource allocation but also empowers the organization with a sustainable, scalable recommendation solution for future developments.

### 8.4 Caveats and limitations

Data drift can cause a significant decrease in performance, requiring a model re-training. When new programmes are added to the catalogue, their feature vectors need to be encoded, the encoded vectors need to be transformed in embeddings, and the new cosine similarities need to be re-calculated. Also, both the encodings and the embeddings need to be cached. When the new programmes don’t share some or all of the tags, the encoding phase produces vectors with some or all zeros, indicative of a loss of information.

## 9 Appendices

### References

- [1] Dor Bank, Noam Koenigstein, and Raja Giryes. Autoencoders. *CoRR*, abs/2003.05991, 2020.
- [2] BBC. BBC Ontologies. <https://www.bbc.co.uk/ontologies/>.
- [3] BBC. Bluey - "more like this" tab. <https://www.bbc.co.uk/iplayer/episodes/m000vbrk/bluey?seriesId=more-like-this>.
- [4] BBC. Things. <https://www.bbc.co.uk/things>.
- [5] BBC. Things - About. <https://www.bbc.co.uk/things/about>.
- [6] BBC. Things - API. <https://www.bbc.co.uk/things/api>.
- [7] BBC. The Royal Charter. <https://www.bbc.com/aboutthebbc/governance/charter>, 2017. Valid until 31 December 2027.
- [8] BBC. How metadata will drive content discovery for the bbc online. <https://www.bbc.co.uk/webarchive/https%3A%2F%2Fwww.bbc.co.uk%2Fblogs%2Finternet%2Fentries%2Feachbb071-d471-4d85-ba9d-938c0c800d0b>, 2020. This page was archived on 1st August 2023 and is no longer updated.
- [9] Tomislav Duricic, Dominik Kowald, Emanuel Lacic, and Elisabeth Lex. Beyond-accuracy: A review on diversity, serendipity and fairness in recommender systems based on graph neural networks, 2023.
- [10] Google. Embeddings. <https://developers.google.com/machine-learning/crash-course/embeddings>.
- [11] Google. Measuring similarity from embeddings. <https://developers.google.com/machine-learning/clustering/dnn-clustering/supervised-similarity>.
- [12] UK Government. UK GDPR. <https://www.legislation.gov.uk/eur/2016/679/contents>, 2016. Regulation (EU) 2016/679 of the European Parliament and of the Council.
- [13] Marius Kaminskis and Derek G. Bridge. Diversity, serendipity, novelty, and coverage. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 7:1 – 42, 2016.

- [14] Lisha Li, Kevin G. Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *CoRR*, abs/1603.06560, 2016.
- [15] Umberto Michelucci. An introduction to autoencoders. *CoRR*, abs/2201.03898, 2022.
- [16] pandas via NumFOCUS, Inc. MultiIndex / advanced indexing. [https://pandas.pydata.org/docs/user\\_guide/advanced.html](https://pandas.pydata.org/docs/user_guide/advanced.html), 2024.
- [17] Simone Spaccarotella. Recommendation assumptions. <https://www.slideshare.net/slideshow/recommendations-assumptions/236291920>, 2015. Presented at Prototyping Day @ Mozilla London on 3 September 2015, at Engineering Summit @ BBC on 7 March 2018, uploaded on SlideShare on 27 June 2020.
- [18] W3C. RDF 1.1 Concepts and Abstract Syntax. <https://www.w3.org/TR/rdf11-concepts>, 2014. W3C Recommendation 25 February 2014.
- [19] W3C. RDF 1.1 Turtle - Terse RDF Triple Language. <https://www.w3.org/TR/turtle/>, 2014. W3C Recommendation 25 February 2014.
- [20] W3C. Resource Description Framework (RDF). <https://www.w3.org/RDF>, 2014. Publication date: 2014-02-25 (with a previous version published at: 2004-02-10).