# Data Science in e-commerce (Search)
## (Talk by : ManojKumar Kannadasan)

## Introduction

In today's digital age, e-commerce platforms have transformed the way we shop, offering unparalleled convenience, variety, and accessibility. At the heart of this transformation lies the power of data science, which plays a pivotal role in enhancing user experiences, optimizing search functionalities, and driving business strategies. Among the giants in the e-commerce sector, eBay stands out not only for its vast marketplace but also for its innovative use of data science, especially in search optimization.

This report delves deep into the multifaceted applications of data science within the e-commerce search optimization domain, with a specific focus on eBay's practices. We explore the journey of data from its acquisition to its final application, discussing challenges, methodologies, and the profound impacts on both sellers and buyers. Through topics such as query categorization, personalized query autocompletion, and spell correction, we'll uncover the intricacies of ensuring that users find exactly what they're looking for, seamlessly and efficiently.

As we navigate through the nuances of these applications, readers will gain insights into the dynamic interplay between data, algorithms, and user experience, highlighting the transformative power of data science in e-commerce.

**The base questions :**
**1. Describe the market sector or sub-space covered in this lecture.**

The market sector or sub-space covered in this lecture pertains to e-commerce, specifically focusing on search optimization and user experience improvements within an e-commerce platform. The lecture delves deep into the data science applications at eBay, one of the world's largest e-commerce platforms. Topics such as search functionalities (text, image, voice, and conversational search), query categorization, and personalized search suggestions are highlighted. These features and functionalities are critical in ensuring users can find and purchase products efficiently and effectively on such platforms. Furthermore, the lecture emphasizes the importance of data-driven decision-making in enhancing the search experience, correcting user input errors (like misspellings), and providing personalized user experiences.

**2. What data science related skills and technologies are commonly used in this sector?**
**Based on the lecture content, the following data science-related skills and technologies are commonly used in the e-commerce search optimization sector:**

1. Data Acquisition, Extraction and Mining: Skills in gathering data from various sources, which could include user interactions, transaction records, click streams, etc. The ability to uncover patterns and insights from vast amounts of data.

2. Clustering/Classification: Techniques to group similar data points or categorize them into predefined classes.

3. Data Modeling and Exploratory Data Analysis: Building predictive or descriptive models based on historical data. The initial phase of understanding and visualizing the data to find patterns or anomalies.:

4. Data Warehousing: Storing, retrieving, and managing large datasets efficiently.

5. Predictive Analytics and Regression Analysis: Using statistical algorithms and machine learning techniques to identify the likelihood of future outcomes. Understanding relationships between variables.

6. Data Visualization: Tools and techniques to visually represent data, aiding in understanding and decision-making.

7. Machine Learning: Especially supervised learning for tasks like query categorization, and unsupervised learning for tasks like clustering.

8. Deep Learning: Using neural networks, such as Convolutional Neural Networks (CNNs), as mentioned in the Convolutional Latent Semantic Model (CLSM) for eBay's Query Categorization.

9. Text Mining & Natural Language Processing (NLP): Extracting useful information from unstructured text data. For tasks like spell correction, query autocompletion, and understanding user intent in searches.

10. Search Algorithms: For optimizing search results, ranking, and relevancy.

11. Recommendation Systems: Algorithms that suggest products to users based on their behavior and preferences.

12. Efficiency Optimization: Techniques to ensure that search and recommendation systems are fast and responsive, even with massive amounts of data.

13. Error Modeling: Understanding and predicting the types of errors users might make, such as typos or misspellings, and offering corrections.

14. Word Embeddings: Like Word2Vec, fastText, and GloVe for converting words into vectors, aiding in tasks like similarity checks and contextual understanding.

15. Learning to Rank (L2R): Machine learning algorithms specifically designed for ranking problems.

The e-commerce search optimization sector demands a blend of these skills and technologies to ensure users have a seamless, efficient, and personalized shopping experience.

**3. How are data and computing related methods used in typical workflows in this sector? Illustrate with an example.**

In the e-commerce search optimization sector, data and computing-related methods play a pivotal role in enhancing user experience, improving search accuracy, and ensuring efficient product discovery. The integration of data science and computational methods transforms raw data into actionable insights and user-friendly features. Here's a breakdown of a typical workflow, illustrated with the example of "Personalized Query Autocompletion" from the provided lecture:

Workflow:
1. Data Collection: Begin by gathering vast amounts of data. This can be user behavioral data, search queries, transaction history, and more. For personalized query autocompletion, data is collected from billions of user sessions. This includes prefixes typed by the user, queries clicked from autocomplete suggestions, previous queries issued by the user, and global performance metrics of the query.

2. Data Pre-processing: Clean and preprocess the data to remove anomalies or irrelevant information. For the query autocompletion, the user's search history in a session is analyzed to predict their intent. Data preprocessing ensures that only relevant and clean search histories are used for predictions.

3. Feature Engineering: Transform the data into a format suitable for machine learning models. This involves creating meaningful attributes or features from the raw data. In our example, features are computed based on previous queries issued by the user. Textual features like n-grams, frequency, session-based metrics, and similarity scores (based on text and vector representations) are generated.

4. Model Training: Using the processed data, train a machine learning or deep learning model to perform the desired task. In the context of personalized query autocompletion, a Machine Learned Ranking Model is trained. Positive samples (queries clicked in autocomplete) and negative samples (queries viewed but not clicked in autocomplete) are used for training.

5. Evaluation: Once the model is trained, it is evaluated using various metrics to ensure its accuracy and effectiveness. For the query autocompletion, metrics like MRR (Mean Reciprocal Rank), Success Rate, MAP (Mean Average Precision), and nDCG (Normalized Discounted Cumulative Gain) are used.

6. Deployment: After ensuring the model's effectiveness, it is deployed in the real-world environment. In the case of eBay, once the personalized query autocompletion model is trained and evaluated, it is integrated into their search platform. As users start typing, the model provides real-time personalized suggestions based on their search history and behavior.

7. Feedback Loop: Continuously monitor the model's performance in the real-world setting. Gather feedback and data to further refine and retrain the model. For instance, if users consistently ignore or avoid certain autocomplete suggestions, this feedback can be used to fine-tune the model.

**Illustration:**
Imagine a user who frequently searches for camera equipment on eBay. They've searched for terms like "DSLR camera," "Canon EOS," and "camera lens" in the past. Now, when they begin typing "Can...", instead of showing generic popular suggestions like "candle" or "canvas," the personalized query autocompletion system, trained on their past behavior, might suggest "Canon EOS" or "Canon camera lens," making their search experience faster and more relevant.

This workflow, powered by data and computational methods, ensures that e-commerce platforms like eBay can provide a tailored and efficient search experience to their users.

**4. What are the data science related challenges one might encounter in this domain?**
**The e-commerce search optimization sector, while ripe with opportunities for data science applications, also presents a myriad of challenges. Here are some data science-related challenges one might encounter in this domain:**

1. Volume of Data: E-commerce platforms often deal with massive amounts of data due to the sheer number of transactions, user interactions, and product listings. Managing, processing, and drawing insights from such vast data sets can be challenging.

2. Data Quality: The data collected might be noisy, incomplete, or even erroneous. For instance, user-generated content (like product reviews or seller descriptions) can be inconsistent and unstructured.

3. Real-time Processing: Features like query autocompletion or personalized recommendations often require real-time processing to provide instant feedback to users. Achieving this in the presence of vast data and without compromising accuracy is challenging.

4. Scalability: As the platform grows, the solutions need to scale efficiently. A model or infrastructure that works for a million users might not work as efficiently for a hundred million.

5. Diverse User Behavior: Users from different regions, backgrounds, or cultures might have varying behaviors and preferences. Building a one-size-fits-all model can be ineffective.

6. Changing Trends: E-commerce is influenced by seasonal trends, new product launches, or even current events. The models need to be adaptive to these changing trends.

7. Cold Start Problem: For new products or new users, there's limited data to base recommendations or predictions on. Addressing this cold start problem is a common challenge.

8. Bias and Fairness: Models might inadvertently introduce or amplify biases present in the training data, leading to unfair or biased recommendations.

9. Interpreting Complex Models: While deep learning models might provide better accuracy, they are often seen as black boxes. Interpreting their predictions in meaningful ways, especially in business contexts, can be challenging.

10. Ensuring Privacy: With increasing concerns about user privacy, using user data for personalization while ensuring data privacy and compliance with regulations (like GDPR) is a balancing act.

11. Spell Corrections & Language Nuances: Handling typos, regional language variations, slang, or new internet lingo can be tricky but is essential for search optimization.

12. Multimodal Data: E-commerce platforms are increasingly incorporating image and voice search. Handling and deriving insights from such multimodal data adds another layer of complexity.

13. Feedback Loops: If not managed correctly, feedback loops can arise, where the system reinforces its own predictions. For instance, if a product is recommended frequently and thus clicked on more, the system might view it as increasingly relevant and recommend it even more.

14. Over-Personalization: While personalization enhances user experience, overdoing it can limit the user's exposure to a variety of products, potentially affecting sales and user discovery of new items.

15. Evaluation Metrics: Determining the right evaluation metrics that align with business goals is essential. Traditional accuracy might not always reflect the real-world effectiveness of a model.

Addressing these challenges requires a combination of technical expertise, business acumen, and continuous iteration and learning. However, overcoming them can lead to significant improvements in user experience and business outcomes.

**5. What do you find interesting about the nature of data science opportunities in this domain?**
**The nature of data science opportunities in the e-commerce search optimization domain is fascinating for several reasons:**

1. Multifaceted Impact: Data science in e-commerce can influence various aspects – from improving user experience with personalized recommendations to optimizing supply chain logistics. The breadth of impact is vast.

2. Real-time Decision Making: The need for real-time decision-making, such as instant product recommendations or query autocompletions, presents both a challenge and an opportunity to design efficient algorithms that work at scale.

3. Multimodal Data Integration: E-commerce platforms are increasingly embracing diverse data types – text, images, voice, and even augmented reality. Integrating insights from these varied data sources presents a rich and complex opportunity.

4. Dynamic Nature: E-commerce is influenced by rapidly changing trends, be it seasonal changes, new product launches, or viral internet trends. This dynamic nature requires models that can adapt quickly, making the problem-solving aspect very engaging.

5. Global User Base: E-commerce platforms cater to a global audience. Understanding and catering to the diverse cultural, linguistic, and behavioral nuances of this vast user base is an exciting challenge.

6. Feedback-driven Iteration: The direct feedback loop, where changes can be immediately tested and their impacts measured (e.g., A/B testing), allows for rapid iteration and improvement. This iterative nature is exciting for anyone who loves to continually refine and optimize.

7. Ethical Considerations: With the power of personalization comes the responsibility of ensuring user privacy and fairness. Navigating these ethical considerations adds depth to the technical challenges.

8. Economic Impact: Optimizations in this domain can lead to substantial economic impacts. A slight improvement in recommendation accuracy or search relevance can translate to significant revenue increases for large platforms.

9. Interdisciplinary Collaboration: The domain necessitates collaboration between data scientists, UX designers, business strategists, and engineers. This interdisciplinary nature ensures a holistic approach to problem-solving.

10. Innovation in Algorithms: The sheer volume of data and the need for efficiency drive innovation in algorithms. For instance, the evolution from traditional recommendation systems to deep learning-based ones or the integration of NLP for better search are testaments to this innovation.

11. User-Centric Focus: At the heart of all these opportunities is the end-user. Ensuring a seamless, personalized, and efficient user experience makes the work in this domain truly impactful.

In essence, the e-commerce search optimization domain offers a unique blend of technical challenges, rapid innovation, and direct, measurable impact. It's a domain where data science can tangibly enhance millions of daily online shopping experiences, making it both rewarding and intriguing.


**(i). Please discuss how sellers and buyers may need different data features in an e-commerce platform such as e-Bay. (10 pts of the 80 C+R points in the rubric)**

In an e-commerce platform like eBay, both sellers and buyers interact with the platform, but their needs, objectives, and interactions differ significantly. As a result, the data features they require are distinctively tailored to their specific roles and goals. Let's explore the different data features needed by each group:

**Sellers:**
1. **Listing Analytics:** Detailed statistics about how their product listings are performing in terms of views, clicks, and sales.

2. **Sales Metrics**: Historical and real-time sales data, average transaction values, return rates, and overall revenue.

3. **Inventory Management** : Tools to monitor stock levels, predict inventory needs, and optimize the timing and quantity of restocks.

4. **Pricing Insights** :Data about competitor pricing, historical pricing trends, and price elasticity to help in setting competitive prices.

5. **Buyer Behavior Analytics** : Insights into what buyers are looking for, commonly viewed products, and common reasons for cart abandonment.

6. **Feedback & Reviews Analytics**: Aggregated data on customer feedback and reviews, with the ability to drill down into specific comments or issues.

7. **Shipping & Logistics**: Data on shipping durations, costs, and feedback to optimize shipping methods and providers.

8. **Ad Performance Metrics**: For sellers who advertise on the platform, metrics on ad impressions, clicks, conversion rates, and return on advertising spend.

9. **Return & Dispute Metrics**: Insights into return rates, common reasons for returns, and dispute resolutions.

10. **Market Trends**: Data on emerging market trends, popular product categories, and seasonal demand changes.

**Buyers:**
1. **Personalized Recommendations**: Suggestions based on browsing history, past purchases, and preferences.

2. **Search Features**: Enhanced search capabilities, including text, image, and voice search, with real-time autocompletion and correction.

3. **Product Reviews & Ratings**: Aggregated ratings and detailed reviews from other buyers to inform purchasing decisions.

4. **Price Comparison Tools**: Features to compare prices of a product across different sellers or similar products.

5. **Wish Lists & Tracking**: Ability to create wish lists, track price changes, and receive notifications for desired products.

6. **Purchase History & Analytics**: Detailed records of past purchases, spending analytics, and insights into spending habits.

7. **Shipping & Delivery Information**: Real-time tracking of shipments, historical data on delivery times, and options to choose preferred delivery methods.

8. **Discounts & Offers**: Information on available discounts, loyalty points, and personalized offers.

9. **Safety & Trust Metrics**: Data on seller trustworthiness, based on ratings, reviews, and dispute resolutions.

10. **Customer Support Interactions**: History of interactions with customer support, resolution rates, and feedback options.

The sellers are primarily focused on optimizing their sales, understanding their market position, and improving their offerings, buyers are more focused on finding the right products at the best prices and ensuring a seamless purchasing experience. The data features on an e-commerce platform like eBay are thus tailored to cater to these unique needs and objectives of both groups.

**(ii) Describe briefly the algorithmic steps involved in query correction as described in the lecture. (10 pts of the 80 C+R points in the rubric)**

The process of query correction, as described in the lecture, aims to correct potential misspellings or inaccuracies in user search queries to improve search results. The algorithmic steps involved are:

1. **Candidate Generation (Efficiency):** For a given incorrect query, generate a set of possible correct query candidates. This can involve looking at close matches in terms of spelling, phonetics, or keyboard proximity. Techniques like tries (a tree-like data structure) can be employed to efficiently generate and store these candidates.

2. **Language Model (Big & Special)**: This model evaluates how likely a candidate correction is based on its occurrence in the data or its syntactic and semantic correctness. For instance, a query like "levis blue jeans 32 in" would be evaluated based on its occurrence frequency. If a particular phrase or word is rarely seen or doesn't make linguistic sense, it's less likely to be a valid correction. The model often employs the Markov assumption (especially second-order) to predict the likelihood of a sequence of words.

3. **Error Model (Big & Special)**: This model assesses how likely it is that a user intended to type the correct query but ended up typing the incorrect one. It considers factors like keyboard distance (how close keys are on a keyboard, leading to typos) and phonetic distance (how words sound, leading to phonetic misspellings). The model is trained using triples of intended words, observed words, and their counts. For example, the likelihood of "the" being typed as "teh" is high, whereas

"the" being typed as "hippopotamus" is extremely low. This model can be constructed by analyzing logs to see common misspellings and their intended corrections.

4. **Ranking (Precision):** After generating candidates and evaluating them using the language and error models, rank the corrections based on their combined likelihood. This involves a mathematical formulation where, for each candidate correction for a word, you compute times , where is the probability of the correction (from the language model) and is the probability of observing the incorrect word given the correction (from the error model). The candidate with the highest combined probability is selected as the best correction.

5. **Output Corrected Query**: Present the top-ranked corrected query to the user or use it to fetch and display search results.

The query correction process involves generating potential corrections, evaluating their linguistic likelihood, assessing the likelihood of typing errors, and then combining these factors to rank and select the best possible correction.

**(iii) Also, answer the following multiple-choice questions: You can list the quetsion number and the letter corresponding to the correct choice as Answer in your report, (2x5 = 10 pts of the 80 C+R points in the rubric)**

**Q1 C**
**Q2 B**
**Q3 C**
**Q4 B**
**Q5 D**

**Reference:**
[1] Lecture video Lecture 5: ManojKumar Kannadasan
[2]
[3] Article - https://en.wikipedia.org/wiki/Learning to rank
[4]
[5]