# Overall Data Science
# Review Topic List

This list is a high-level overview of all things related to data science that are important to have an understanding for in order to (1) feel confident in data science skills, (2) be successful in data science job hunting & (3) be successful in a data science job. Not every job will make use of every one of these concepts, but all have potential to be so, or be tested during an interview process.

Legend

| Abbreviation | Description | Notes |
|---|---|---|
| FK (for topic X, Y, Z) | Foundational Knowledge - topic/skill that provides foundations for understanding more complex topics or carrying out other responsibilities as a data scientist | Parentheticals describe which more complex topics the FK is required for |
| DSIK | Data Science Interview (Knowledge) - commonly tested topic during conversational portions of data science interviews | Topics can be addressed in several formats:<br>- Q & A (conversational "quizzing")<br>- Whiteboarding (pseudoalgorithm) |
| DSICC | Data Science Interview (Code Challenge) - commonly tested topic addressed in code challenges | Code challenges can come in the forms of:<br>- Take home<br>- Live coding (on laptop, whiteboard, etc.) - live coding usually involves less complex challenges (basic SQL, python coding, etc.) |
| NTH | Nice To Have - Not absolutely required, but will be helpful to know, more likely to get a job | Can help data scientists stand out during interview process, and will help facilitate more in-depth learning & greater complexity of development |
| OPN | Optional - Completely optional | Can help data scientists stand out during interview process |

| Topic | Subtopics | Overall Priority | Relevance (legend) | Notes |
|---|---|---|---|---|
| **Coding** | Python<br>- Basics<br>- Data structures<br>- OOP<br>- Libraries | Very High | - FK (for most of DS implementation)<br>- DSICC<br>- DSIK (sometimes) | - Sometimes topics can be tested like OOP, BigO, multi threading, etc. |
| | SQL<br>- Order of operations<br>- Joins<br>- MySQL<br>- CloudSQL | Very High | - FK (for data wrangling & management)<br>- DSICC<br>- DSIK | - Sometimes topics can be tested verbally like SQL order of operations, join vs union, etc. |
| | R | Low | - OPN (assuming python) | - Becoming optional in DS |
| | Bash / command line | Medium | - NTH | - 'Required' for git use for the most part<br>- Especially useful for interacting with cloud technology |
| | Git<br>- Push<br>- Pull<br>- Rebase<br>- Forking | High | - FK (for collaborative work streams)<br>- NTH | - Not often tested during interviews, but may be required to download data for code challenge<br>- Will likely be used in day-to-da work in DS job |

| Topic | Subtopics | Overall Priority | Relevance (legend) | Notes |
|---|---|---|---|---|
| **Statistics** | Basic estimates<br>- Mean<br>- Median<br>- Mode<br>- Percentiles<br>- Boxplots<br>- Variability vs Standard deviation | High | - FK (for cohort analytics)<br>- DSIK<br>- DSICC | Can be tested both on a knowledge basis, or expected during code challenge |
| | Probability<br>- Calculating probability with different scenarios<br>- Probability & cumulative density functions, kernels | High | - FK (for segmentation, clustering, hyperparameter optimization, a/b testing)<br>- DSIKCC | |
| | Hypothesis testing<br>- H0, Ha<br>- Meaning & interpretation of p-value<br>- Significant testing (t-test, z-test, chi-square)<br>- Confidence intervals<br>- Bootstrapping / resampling<br>- Sample sizes / power analysis<br>- Effect size<br>- A/B testing<br>- Multi arm a bandit | Very High | - FK (for A/B testing, academic research with health data, marketing strategy decision making, for hyperparameter optimization)<br>- DSIK<br>- DSICC | |
| | Concepts / rules<br>- Types of bias<br>- Central limit theorem<br>- Sample vs Population estimates<br>- Bootstrapping<br>- Permutation test | Very High | - FK (for building models)<br>- DSIK | |

| Topic | Subtopics | Overall Priority | Relevance (legend) | Notes |
|---|---|---|---|---|
| | Distributions<br>- Normal / Gaussian<br>- Binomial<br>- Poison<br>- Exponential<br>- Plotting<br>- Using for: A/B testing, assumption-testing | Very High | - FK (for data preparation & determining validity of prediction results)<br>- DSIK | Need to know assumptions / description of each, what they're used for, etc. |
| **Data Processing** | - Importing CSV, JSON, txt, etc<br>- Categorical & dummy variables<br>- Time stamps & date | Very High | - FK (for most of data science implementations)<br>- DSICC | - Many different libraries & standards for time stamps and they interact differently with various plotting libraries |
| **Data Exploration** | Topics<br>- Examining data (nulls, completeness, etc.)<br>- Imputing data<br>- Data point density | Very High | - FK (for examining feature relationships, discovering bias, time series data that is incomplete)<br>- DSICC | - Imputing data mostly relates to time series data |
| | Python coding<br>- Matplotlib<br>- Seaborn<br>- Plotly | Very High | - FK<br>- DSICC | |
| | Approaches<br>- Feature comparison by plotting<br>- Regression<br>- Scatterplots<br>- Histograms<br>- Sub-segmentation | Very High | - FK<br>- DSICC | |

| Topic | Subtopics | Overall Priority | Relevance (legend) | Notes |
|---|---|---|---|---|
| **Modeling** | Concepts<br>- Simplicity vs need for complexity<br>- Bias vs variance<br>- Supervised vs unsupervised vs reinforcement vs semi supervised<br>- How to select model type | Very High | - FK (for determining MVP = minimum viable product, resource requirements with DevOps)<br>- DSIK | - How bias vs variance trade off affects model parameter tuning & structure selection<br>- Over concepts will be HEAVILY tested in knowledge-based interview stages |
| **(Supervised)** | Regression<br>- Assumptions<br>- Least squares linear regression<br>- Understanding of multiple linear regression<br>- Dummy variables<br>- Testing / scoring<br>- Regularization (L1 & L2)<br>- Polynomial / Spline, etc. | Very High | - FK (feature engineering & dimensionality reduction)<br>- DSICC<br>- DSIK<br>- Some NTH | - WIll be tested heavily in interviews both conceptually (bias/variance, scoring methods, assumptions) & coding<br>- Advanced approaches (polynomial / splines are NTH) |
| **(Supervised - Regression)** | Evaluation & scoring<br>- Residuals (error): absolute error, squared error, MSE, RMSE, A1C, B1C<br>- Parameter selection (Regularization) | Very High | - FK  (for model selection & hyperparameter optimization)<br>- DSIK | - WIll be tested heavily in interviews both conceptually |
| **(Supervised)** | Classification<br>- Class balancing<br>- Naive Bayes (using probability)<br>- Logistic regression<br>- Gradient Boosted Classifier<br>- Odds ratio | Very High | - FK (feature engineering, clustering)<br>- DSIK<br>- DSICC | - Foundational knowledge that is very important, and also commonly required for implementation in take home assignments or code challenges |

| Topic | Subtopics | Overall Priority | Relevance (legend) | Notes |
|---|---|---|---|---|
| **(Supervised - classification)** | Classification evaluation<br>- Confusion matrix<br>- Precision<br>- Recall (sensitivity)<br>- Specificity<br>- AUC<br>- Lift | Very High | - FK<br>- DSIK<br>- DSICC | - Extremely high priority for DS interview - knowledge testing!! Especially Precision, recall, specificity, AUC |
| **(Supervised)** | Decision Trees<br>- Single vs Ensemble<br>- Purity (Gini, entropy, etc.)<br>- Information Gain (split decision)<br>- Parameters (leaf depth, randomization, etc.)<br>- Bagging<br>- Random Forest<br>- Boosting | Very High | - FK (feature engineer, clustering)<br>- DSIK<br>- DSICC | - Foundational knowledge that is very important, and also commonly required for implementation in take home assignments or code challenges |
| **(Unsupervised)** | Clustering<br>- Standardization / PCA beforehand<br>- Distance metrics<br>- N-cluster selection<br>- KNN<br>- Kmeans<br>- SVD<br>- Use cases, advantages, disadvantages<br>- Hierarchical | Very High | - FK (for recommendation engines, feature engineer)<br>- DSIK | - Foundational knowledge that is very important, and also commonly required for implementation in take home assignments or code challenges |

| Topic | Subtopics | Overall Priority | Relevance (legend) | Notes |
|---|---|---|---|---|
| (Supervised/ Semisupervised) | Neural Networks<br>- Forward propagation<br>- Backward propagation<br>- Hidden layer<br>- Activation functions (tanh, sigmoid, ReLu)<br>- Types (CNN, RNN, LSTM, Autoencoder, transformers, GAN, +++)<br>- Tensorfow / Keras | High - especially for MLE | - FK<br>- DSIK | - Usually tested conceptually. For MLE interviews, you will definitely need to be intimately familiar with deep learning at large, including these topics, know how to code, deep dive into parameter discussion, etc. |
| Model Selection | Model Comparison<br>- Simplicity vs complexity<br>- Bias vs Variance<br>- Data completeness & availability<br>- Goal<br>- Reproducibility<br>- Personalization<br>- Model Versioning | High | - FK (for selecting MVP, for reproducibility improvement, for versioning)<br>- DSIK | - Concepts tested to help interviewers know what your approach to problem solving with DS is |
| | Model outcome interpretation<br>- Over/under-fitting<br>- Train vs validation vs Test scores<br>- Feature importance<br>- Hyperparameter optimization | Very High | - FK<br>- DSIK | - WIll be tested heavily (concepts) during conversational parts of DS interviews |
| Text data | Natural Language Processing (NLP)<br>- tf-idf<br>- Stop words<br>- Tokens | High/medium | - FK<br>- DSIK<br>- NTH | - May be asked about familiarity, especially if it's relevant to the job |

| Topic | Subtopics | Overall Priority | Relevance (legend) | Notes |
|---|---|---|---|---|
| **Data Tooling** | Big Data Tools<br>- Cloud platforms (AWS, GCP, Azure)<br>- Spark<br>- HDFS (old, but still asked about)<br>- Dask<br>- Multi-threading (python) | Medium | - FK<br>- NTH | - May be asked about familiarity, especially if it's relevant to the job |
| | Spark<br>- Resilient Distributed Dataset<br>- Lazy evaluation<br>- MapReduce<br>- Spark SQL<br>- PySpark<br>- Spark Mlib | Medium | - FK<br>- NTH | - Differences between Spark & HDFS/MapReduce are common test questions |
| | Databases<br>- Centralized vs distributed<br>- On-prem vs Cloud<br>- Relational vs Non-relational (NoSQL, document-based, etc.) | High | - NTH<br>- DSIK | - May be asked about familiarity, especially if it's relevant to the job |
| | Cloud Tools<br>- Databases<br>- Virtual machines<br>- Service Accounts<br>- Ingestion pipelines (e.g. AWS Kinesis/firehose, GCP Dataflow) | Medium/High | - NTH | - May be asked about familiarity, especially if it's relevant to the job |
| | Other Tools<br>- API / django/ flask<br>- Pub/sub (or similar)<br>- Tensorflow / Keras | Low | - NTH<br>- OPN | - May be asked about familiarity, especially if it's relevant to the job |
| | | | | |