

Scala Coding Guidelines

Quantexa Coding Guidelines:

General Scala/Spark	
Camel Case	<p>Use camel casing, not underscores! Vals/vars/methods(defs) should all be lowerCamelCase, types/objects should all be UpperCamelCase (also known as ProperCase), for further details please see this style guide.</p> <p>Exceptions</p> <ol style="list-style-type: none">1. In ETL processes, the raw data should be read into a parquet maintaining all of the source field names exactly.2. Constant values - use ProperCase. For example, val YearInDays = 3653. When using acronyms and camel case, if it is the first word within a variable name it should be lower case, if it is not the first word, it should be all upper case. E.G:<ul style="list-style-type: none">○ jiraTaskDescription (and not JIRATaskDescription)○ accountIBAN (and not accountIban)4. When using acronyms and proper case, every letter of the acronym should be upper case, e.g.<ul style="list-style-type: none">○ case class ETLConfig(hdfsPath: String)5. Packages all lower case, including package objects.
Abbreviation (don't do it!)	<p>Dont abbreviate words in names, for example:</p> <ul style="list-style-type: none">• accountDate (not accountDt)• accountDateTime (not accountDtTm, and not accountDatetime) <p>Exceptions</p> <ul style="list-style-type: none">• When using Id (for Identification or Identifier), do not shorten and use as lowerCamelCase, for example accountId
Split long lines of code	<p>Long lines of code should be split over multiple lines. The Scala style guide suggests 80 characters as an upper bound, but discretion is allowed (R&D use 120 because they have big monitors!). If it doesn't fit on your screen without scrolling then the line definitely needs splitting!</p>

General Scala/Spark	
Comments	<p>If you need to comment your code to explain something, instead question whether you can write the code in a clearer way (generally the way to do this should result in more functional code). Good code should read like a story; if more detail is required, the developer can click through to see more information. A simple example code which follows this notion is as follows:</p> <p>Examples of where comments are normally required:</p> <ol style="list-style-type: none"> 1. Where decisions are being made which can't be explained through code (e.g. features of the data and business feedback which determined how to write a score) 2. "TODO"s and "FIXME" comments (see below)
Scaladocs	<p>See style guide. Scaladocs are effectively comments about what something does and for methods/functions the parameters it takes. Objects and key methods/functions should have Scaladocs explaining what they are/do. An example from the style guide:</p> <pre> 1. /** Creates a person with a given name and birthdate 2. * 3. * @param name their name 4. * @param birthDate the person's birthdate 5. * @return a new Person instance with the age determined by the 6. * birthdate and current date. 7. */ 8. def apply(name: String, birthDate: java.util.Date) = {} 9. }</pre>

General Scala/Spark	
TODO/FIXME	<p>When something is left to do, use a TODO:</p> <pre>//TODO: Add support for Longs</pre> <p>When something may not work as expected, use a FIXME:</p> <pre>//FIXME: Pattern match will fail if user provides Integer parameters</pre> <p>In both cases, be specific about what the task is, as it may be a different developer who works on the task.</p>
Hard-coded configuration	<p>Do not hard code configuration! Examples which should be externalised into a configuration project (e.g. type-safe) include:</p> <ol style="list-style-type: none"> 1. Locations of raw data / output data 2. Iteration number 3. Dates 4. Scoring Parameters where the client wishes to be able to change these parameters
Regular Expressions	<p>Regular expressions are compiled at run time, and add significant overhead. As such, where possible they should be defined once only. In other words:</p>

Spark only	
Broadcast joins	<p>When one dataset in a join is small it is a good idea to suggest to Spark to broadcast the small dataset.</p> <pre>largeDF.join(broadcast(smallDF), Seq("foo"))</pre> <p>A broadcast join is sort of like a hash merge in SAS - instead of shuffling both datasets to do a join the smaller one is 'sent' to each executor which means the large dataset doesn't need to 'move'. This is more efficient.</p>

Joins on same column(s)	<p>Use <code>df1.join(df2, Seq("id"))</code></p> <p>Dont Use <code>df1.join(df2, df1("id") === df2("id"))</code></p> <p>This is only applicable where you are joining two tables on a column (or columns) with the same name. This avoids ending up with duplicate columns in the output dataframe and is much more succinct, especially when joining on multiple columns</p>
Prefer select over drop	<p>Instead of doing a long series of drops</p> <p>e.g. <code>df.drop("a").drop("b").drop("c").drop("d").drop("e")</code></p> <p>consider selecting instead</p> <p>e.g. <code>df.select("g","h","i","j","k")</code></p> <p>Selecting makes it clear what is on the output dataframe. Furthermore, Spark will not complain if you try and drop a column which doesn't exist (which in general is not desirable).</p> <p>Additional benefits of select over drop include being able to choose the order of the columns and also rename them in the same step: <code>df.select(\$"g".as("moreLogicalName"),...)</code></p>
Column referencing	<p>There are 3 ways of referencing columns:</p> <ul style="list-style-type: none"> • <code>col("columnName")</code> • <code> \$"columnName"</code> • <code>'columnName'</code> <p><code> \$"columnName"</code> is preferred as it is short, but won't ruin code highlighting if you have code in a text editor such as notepad++. For this reason, avoid <code>'columnName'</code></p> <p>Note that it is very common to do some coding in notepad++ on projects when doing ad-hoc coding work before moving it into an etl/scoring project.</p>

External comprehensive coding guidelines:

<http://docs.scala-lang.org/style/>

<https://github.com/alexandru/scala-best-practices>

<http://www.lihaoyi.com/post/StrategicScalaStylePracticalTypeSafety.html>

<https://pavelfatin.com/scala-collections-tips-and-tricks/>

<https://github.com/alexandru/scala-best-practices>

https://twitter.github.io/scala_school/