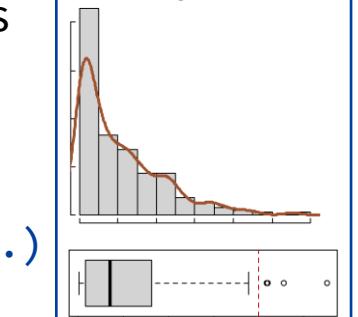
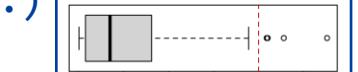
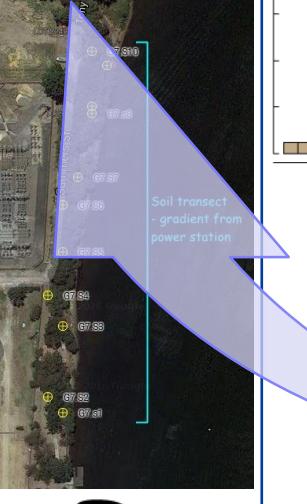
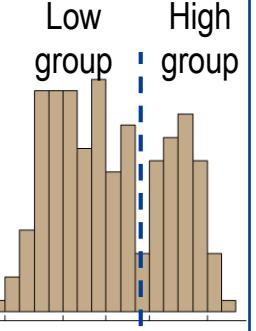
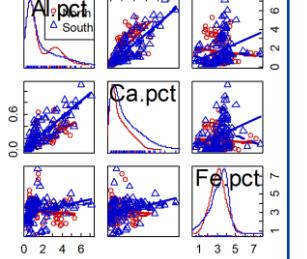
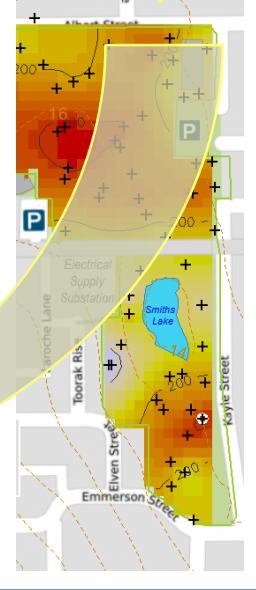


# Basic EDA workflow for environmental assessment

Clean data	Import data	Check variables	New variables	New factors	Relationships	Means comparisons	Spatial patterns																																
<p>Can be done in Excel or R</p> <ul style="list-style-type: none"> <li>• single heading row</li> <li>• no blank rows or columns</li> <li>• censor unreliable values</li> <li>• consistent missing value codes and factor levels</li> <li>• ideally, save with metadata</li> <li>• etc.</li> </ul> 	<p>From csv, Excel, online, shapefile, etc.</p>  <pre>sv2017_export data &lt;- read.csv(...) # and so on...</pre>	<p>Distributions</p> <ul style="list-style-type: none"> <li>• normality</li> <li>• modality</li> <li>• transform if necessary</li> <li>• re-check transformed variables (esp. for modality)</li> </ul> <p>Extremes</p> <ul style="list-style-type: none"> <li>• identify unusual values</li> <li>• identify values over guidelines</li> </ul>  	<ul style="list-style-type: none"> <li>• (including transformed variables)</li> <li>• distance (from upstream, from source, etc.) [use coordinates and Pythagoras]</li> </ul>  	<p>Based on:</p> <ul style="list-style-type: none"> <li>• clearly bimodal or multimodal variables</li> <li>• near/far, etc.</li> </ul> 	<ul style="list-style-type: none"> <li>• bivariate plots e.g. scatter plot matrices</li> <li>• separate by groupings</li> <li>• correlation matrices</li> <li>• outliers in variables</li> </ul> <table border="1"> <caption>Pearson's r</caption> <thead> <tr> <th></th> <th>Al.pct</th> <th>Ca.pct</th> <th>Fe.pct</th> </tr> </thead> <tbody> <tr> <td>Al.pct</td> <td>1.00</td> <td>0.77</td> <td>0.16</td> </tr> <tr> <td>Ca.pct</td> <td>0.77</td> <td>1.00</td> <td>0.18</td> </tr> <tr> <td>Fe.pct</td> <td>0.16</td> <td>0.18</td> <td>1.00</td> </tr> </tbody> </table> <table border="1"> <caption>P</caption> <thead> <tr> <th></th> <th>Al.pct</th> <th>Ca.pct</th> <th>Fe.pct</th> </tr> </thead> <tbody> <tr> <td>Al.pct</td> <td>0.0000</td> <td>0.0232</td> <td>0.0095</td> </tr> <tr> <td>Ca.pct</td> <td>0.0000</td> <td>0.0000</td> <td>0.0000</td> </tr> <tr> <td>Fe.pct</td> <td>0.0232</td> <td>0.0095</td> <td>0.0000</td> </tr> </tbody> </table>		Al.pct	Ca.pct	Fe.pct	Al.pct	1.00	0.77	0.16	Ca.pct	0.77	1.00	0.18	Fe.pct	0.16	0.18	1.00		Al.pct	Ca.pct	Fe.pct	Al.pct	0.0000	0.0232	0.0095	Ca.pct	0.0000	0.0000	0.0000	Fe.pct	0.0232	0.0095	0.0000	<p>Between groups representing:</p> <ul style="list-style-type: none"> <li>• sampling design (e.g. strata)</li> <li>• any new factors</li> </ul> <pre>t.test(...) aov(...) oneway.test(...)  wilcox.test(...) kruskal.test(...)  library(effsize) cohen.d(...)</pre>	<p>Are these:</p> <ul style="list-style-type: none"> <li>• related to sampling design?</li> <li>• related to conceptual site model?</li> <li>• related to any new factors?</li> </ul> 
	Al.pct	Ca.pct	Fe.pct																																				
Al.pct	1.00	0.77	0.16																																				
Ca.pct	0.77	1.00	0.18																																				
Fe.pct	0.16	0.18	1.00																																				
	Al.pct	Ca.pct	Fe.pct																																				
Al.pct	0.0000	0.0232	0.0095																																				
Ca.pct	0.0000	0.0000	0.0000																																				
Fe.pct	0.0232	0.0095	0.0000																																				

EXPLORATORY DATA ANALYSIS