# Machine Learning Algorithm Cheat Sheet

September 09, 2014

Here is a cheat sheet that shows which algorithms perform best at which tasks.

| Algorithm | Pros | Cons | Good at |
|---|---|---|---|
| **Linear regression** | - Very fast (runs in constant time)<br>- Easy to understand the model<br>- Less prone to overfitting | - Unable to model complex relationships<br>-Unable to capture nonlinear relationships without first transforming the inputs | - The first look at a dataset<br>- Numerical data with lots of features |
| **Decision trees** | - Fast<br>- Robust to noise and missing values<br>- Accurate | - Complex trees are hard to interpret<br>- Duplication within the same sub-tree is possible | - Star classification<br>- Medical diagnosis<br>- Credit risk analysis |
| **Neural networks** | - Extremely powerful<br>- Can model even very complex relationships<br>- No need to understand the underlying data<br>– Almost works by "magic" | - Prone to overfitting<br>- Long training time<br>- Requires significant computing power for large datasets<br>- Model is essentially unreadable | - Images<br>- Video<br>- "Human-intelligence" type tasks like driving or flying<br>- Robotics |
| **Support Vector Machines** | - Can model complex, nonlinear relationships<br>- Robust to noise (because they maximize margins) | - Need to select a good kernel function<br>- Model parameters are difficult to interpret<br>- Sometimes numerical stability problems<br>- Requires significant memory and processing power | - Classifying proteins<br>- Text classification<br>- Image classification<br>- Handwriting recognition |
| **K-Nearest Neighbors** | - Simple<br>- Powerful<br>- No training involved ("lazy")<br>- Naturally handles multiclass classification and regression | - Expensive and slow to predict new instances<br>- Must define a meaningful distance function<br>- Performs poorly on high-dimensionality datasets | - Low-dimensional datasets<br>- Computer security: intrusion detection<br>- Fault detection in semiconducter manufacturing<br>- Video content retrieval<br>- Gene expression<br>- Protein-protein interaction |