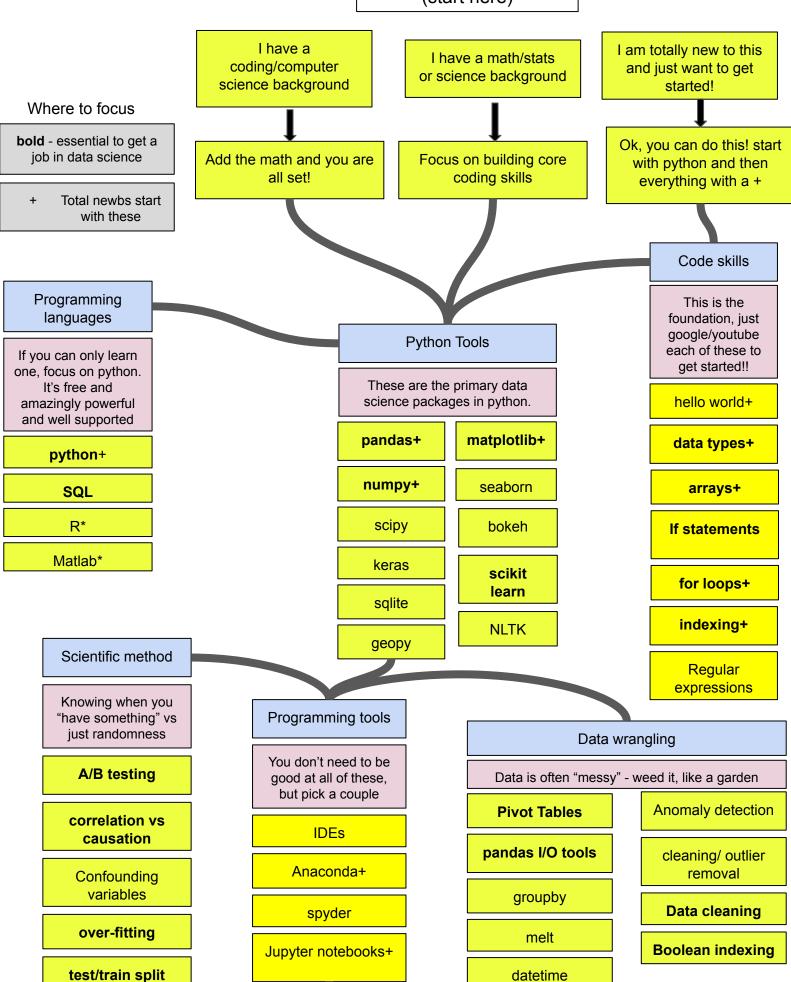
Data Science Roadmap (start here)



Statistics

Dealing with randomness is one of the key jobs of a data scientist. You need to be VERY comfortable with these ideas and tools

Random variables

Regression+

Normal dsn+

Simple linear+

z-scores+

Least squares

t-dsn

Binomial dsn

Multiple linear regression +

Design of experiments

Logistic regression+

Time series

collinearity

Classification

bias - variance tradeoff

Mean, median, correlation, covarinace+

Machine learning skills

Here are the keys to the kingdom! But be careful... don't turn into a brainless button pusher. Make sure you understand the risks in each of these and test early and often.

Hyper-parameter tuning

clustering

Neural nets

k-NN

Classification tress/random forests

k-means+

boosting

unsupervised learning+

Supervised vs

Recommender systems

Gaussian mixture models

Infrastructure

Okay, you have some cool math and ideas of what to do with it... here is what you need to actually turn it loose into production... or "prod" as we call it

github

cron job

Version control

REST API

Front end vs back end

AWS/S3

cloud computing

Metrics

There is no single "best" answer for many problems. These are the ways we measure very complicated things

ROC curve

precision+

AUROC

recall+

specificity

accuracy+

sensitivity

RMSE

Math/Stats skills

Knowing just a *little* bit of theory will save you a ton of time in learning to recognize dead ends BEFORE you spend days coding on a problem

Hypothesis tests

DeMorgan's Rules

independence

Bayes Rule

conditional dsns

Chebyshev's inequality

n choose r

Linear algebra

Almost all of the data you touch will be in some sort of row/column format and learning these tools will unlock some powerful insight for you

Principle Components Analysis

Eigen decomposition

covariance matrix

SVD

Matrix condition number

NMF

Ops Research

Before there was "data science" there was ops research and these tools are essential to solving difficult optimization problems. The good news is that python has a lot of great tools to solve these classic problems.

optimization

Minimax criteria

Nelder-Meed

Queueing theory

Gradient descent

Data meets business

Data science is fundamentally a service discipline, it exists to serve the needs of the organization in which it is practiced. Those organizations are usually businesses, and you need to become versent in the things that businesses care about.

OSEMN process

supply/demand curves

Microeconomics

Price discrimination

CLTV

Price elasticity

nic order

Economic order quantity