

THE ALL-IN- ONE

PYTHON 3,

DATA ANALYTICS,

DATA SCIENCE,

MACHINE LEARNING,

and DEEP LEARNING

CHEAT ~~SHEET~~

BOOK

This Book Was Prepared By:

Jack Reuben Bibi

Machine Learning Engineer

917-353-2666 – jbibi1296@gmail.com

<https://linkedin.com/in/jack-bibi-nyc-> <http://github.com/jbibi1296>

-----SUMMARY-----

I combine curiosity, with passion to build solutions that meet critical challenges. I believe that Machine Learning is a tool that will continue to grow in importance, and I am excited to apply it to important problems.

-----SKILLS-----

PROGRAMMING:

OO Design, Data Structures, Algorithms, Space and Time Complexities, Python, NumPy, Matplotlib, SQL, Scala, AWS, Google Cloud, Jupyter Notebook, APIs, Git, Linux/Unix CLI

DATA:

ETL, Strategy and Planning, Data Mining and Collection, Data Engineering, Data Cleaning, Data Integrity, Data Processing, Data Visualizations and Charts, Excel, Tableau, Communication and Presentation (Technical and Non-Technical)

MACHINE LEARNING:

Keras, TensorFlow, Classification Modeling, Regression Modeling, Supervised Learning, Unsupervised Learning, Ensemble Methods, Neural Networks, NLP, Image Classification, Machine Vision, Time-series, Forecasting

-----TECHNICAL PROJECTS-----

NJ Transit NJCL Track Informer (Python)

I created a customer-centered product that connects with the Twitter API and posts the which track the train will be arriving at. I have received over 20k impressions in the first 48 hours.

Analysis of Colorectal Cancer (Python)

After conducting research, we have created several pipelines in Python to run biostatistical and machine learning analyses to analyze the health of the Colorectal cancer patients.

Salary Predictor (Python)

I designed an app that allows users to upload their resume, choose an algorithm, and get an estimated annual salary. The algorithms were trained on over 6,000 tech and development jobs. This app can is accessible on

<https://salary.jackbibib.dev>

-----EXPERIENCE-----

DEVRY UNIVERSITY – Curriculum Developer

09/2019 – Present

I develop course material and roadmaps for technical courses in the areas of Data Science, Machine Learning, and Software Development. These courses are python-based and use Tensorflow and other Python libraries.

GENIEOUSLY – Data Science Internship

05/2019 – 09/2019

I created forecasting algorithms for clients to assist in Stock Management, Predict Cost, and give an Accurate Pricing Estimate to customers.

OUTDOOR INTL – VP of Marketing and Analytics

06/2017 – 02/2019

I created content, managed, analyzed, and optimized many social media advertisements on Facebook, Instagram, Twitter, and Google SEO.

SOUTHERN TELECOM – Marketing and Analytics Engineer

03/2015 – 01/2017

I researched to find the target audience, target price-point, and desired features of specific product categories to maximize sales. My research increased sales by \$5 million in 12 months by focusing on hot product and taking early advantage of the FAD

-----EDUCATION-----

DEEP LEARNING SPECIALIZATION by Andrew Ng

Coursera 2019

Completed the 5 course specialization in Neural Networks and Deep Learning

ALGORITHMIC TOOLBOX by UC San Diego

Coursera 2019

Data Structures & Algorithms Certificate of Completion 2019

GENERAL ASSEMBLY

02/2019 - 05/2019

Completed a full-time 500+ hour Data Science Immersive program in Python

BERNARD M. BARUCH COLLEGE (CUNY)

BBA Digital Marketing

Data Science Cheatsheet

Compiled by Maverick Lin (<http://mavericklin.com>)
Last Updated August 13, 2018

What is Data Science?

Multi-disciplinary field that brings together concepts from computer science, statistics/machine learning, and data analysis to understand and extract insights from the ever-increasing amounts of data.

Two paradigms of data research.

- Hypothesis-Driven:** Given a problem, what kind of data do we need to help solve it?
- Data-Driven:** Given some data, what interesting problems can be solved with it?

The heart of data science is to always ask questions. Always be curious about the world.

1. What can we learn from this data?
2. What actions can we take once we find whatever it is we are looking for?

Types of Data

Structured: Data that has predefined structures. e.g. tables, spreadsheets, or relational databases.

Unstructured Data: Data with no predefined structure, comes in any size or form, cannot be easily stored in tables. e.g. blobs of text, images, audio

Quantitative Data: Numerical e.g. height, weight

Categorical Data: Data that can be labeled or divided into groups. e.g. race, sex, hair color.

Big Data: Massive datasets, or data that contains greater *variety* arriving in increasing *volumes* and with ever-higher *velocity* (3 Vs). Cannot fit in the memory of a single machine.

Data Sources/Fomats

Most Common Data Formats CSV, XML, SQL, JSON, Protocol Buffers

Data Sources Companies/Proprietary Data, APIs, Government, Academic, Web Scraping/Crawling

Main Types of Problems

Two problems arise repeatedly in data science.

Classification: Assigning something to a discrete set of possibilities. e.g. spam or non-spam, Democrat or Republican, blood type (A, B, AB, O)

Regression: Predicting a numerical value. e.g. someone's income, next year GDP, stock price

Probability Overview

Probability theory provides a framework for reasoning about likelihood of events.

Terminology

Experiment: procedure that yields one of a possible set of outcomes e.g. repeatedly tossing a die or coin

Sample Space S: set of possible outcomes of an experiment e.g. if tossing a die, $S = \{1,2,3,4,5,6\}$

Event E: set of outcomes of an experiment e.g. event that a roll is 5, or the event that sum of 2 rolls is 7

Probability of an Outcome s or P(s): number that satisfies 2 properties

1. for each outcome s , $0 \leq P(s) \leq 1$
2. $\sum p(s) = 1$

Probability of Event E: sum of the probabilities of the outcomes of the experiment: $P(E) = \sum_{s \in E} p(s)$

Random Variable V: numerical function on the outcomes of a probability space

Expected Value of Random Variable V: $E(V) = \sum_{s \in S} p(s) * V(s)$

Independence, Conditional, Compound

Independent Events: A and B are independent iff:

$$P(A \cap B) = P(A)P(B)$$

$$P(A|B) = P(A)$$

$$P(B|A) = P(B)$$

$$\text{Conditional Probability: } P(A|B) = P(A,B)/P(B)$$

$$\text{Bayes Theorem: } P(A|B) = P(B|A)P(A)/P(B)$$

$$\text{Joint Probability: } P(A,B) = P(B|A)P(A)$$

$$\text{Marginal Probability: } P(A)$$

Variability

Standard Deviation Measures the squares differences between the individual elements and the mean

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N-1}}$$

Variance

Interpreting Variance

Variance is an inherent part of the universe. It is impossible to obtain the same results after repeated observations of the same event due to random noise/error. Variance can be explained away by attributing to sampling or measurement errors. Other times, the variance is due to the random fluctuations of the universe.

Correlation Analysis

Correlation coefficients $r(X, Y)$ is a statistic that measures the degree that Y is a function of X and vice versa. Correlation values range from -1 to 1 , where 1 means fully correlated, -1 means negatively-correlated, and 0 means no correlation.

Pearson Coefficient Measures the degree of the relationship between linearly related variables

$$r = \frac{\text{Cov}(X,Y)}{\sigma(X)\sigma(Y)}$$

Spearman Rank Coefficient Computed on ranks and depicts monotonic relationships

Descriptive Statistics

Provides a way of capturing a given data set or sample. There are two main types: **centrality** and **variability** measures.

Centrality

Arithmetic Mean Useful to characterize symmetric distributions without outliers $\mu_X = \frac{1}{n} \sum x$

Geometric Mean Useful for averaging ratios. Always less than arithmetic mean $= \sqrt[n]{a_1 \cdot a_2 \dots a_n}$

Median Exact middle value among a dataset. Useful for skewed distribution or data with outliers.

Mode Most frequent element in a dataset.

Data Cleaning

Data Cleaning is the process of turning raw data into a clean and analyzable data set. "Garbage in, garbage out." Make sure garbage doesn't get put in.

Errors vs. Artifacts

1. **Errors:** information that is lost during acquisition and can never be recovered e.g. power outage, crashed servers
2. **Artifacts:** systematic problems that arise from the data cleaning process. these problems can be corrected but we must first discover them

Data Compatibility

Data compatibility problems arise when merging datasets. Make sure you are comparing "apples to apples" and not "apples to oranges". Main types of conversions/unifications:

- **units** (metric vs. imperial)
- **numbers** (decimals vs. integers),
- **names** (John Smith vs. Smith, John),
- **time/dates** (UNIX vs. UTC vs. GMT),
- **currency** (currency type, inflation-adjusted, dividends)

Data Imputation

Process of dealing with missing values. The proper methods depend on the type of data we are working with. General methods include:

- Drop all records containing missing data
- Heuristic-Based: make a reasonable guess based on knowledge of the underlying domain
 - Mean Value: fill in missing data with the mean
 - Random Value
 - Nearest Neighbor: fill in missing data using similar data points
- Interpolation: use a method like linear regression to predict the value of the missing data

Outlier Detection

Outliers can interfere with analysis and often arise from mistakes during data collection. It makes sense to run a "sanity check".

Miscellaneous

Lowercasing, removing non-alphanumeric, repairing, unidecode, removing unknown characters

Note: When cleaning data, always maintain both the raw data and the cleaned version(s). The raw data should be kept intact and preserved for future use. Any type of data cleaning/analysis should be done on a copy of the raw data.

Feature Engineering

Feature engineering is the process of using domain knowledge to create features or input variables that help machine learning algorithms perform better. Done correctly, it can help increase the predictive power of your models. Feature engineering is more of an art than science. FE is one of the most important steps in creating a good model. As Andrew Ng puts it:

"Coming up with features is difficult, time-consuming, requires expert knowledge. 'Applied machine learning' is basically feature engineering."

Continuous Data

Raw Measures: data that hasn't been transformed yet

Rounding: sometimes precision is noise; round to nearest integer, decimal etc.

Scaling: log, z-score, minmax scale

Imputation: fill in missing values using mean, median, model output, etc..

Binning: transforming numeric features into categorical ones (or binned) e.g. values between 1-10 belong to A, between 10-20 belong to B, etc.

Interactions: interactions between features: e.g. subtraction, addition, multiplication, statistical test distributions more normal), Box-Cox

Statistical: log/power transform (helps turn skewed

max, min, etc

Dimensionality Reduction: using PCA, clustering, factor analysis etc

Discrete Data

Encoding: since some ML algorithms cannot work on categorical data, we need to turn categorical data into numerical data or vectors

Ordinal Values: convert each distinct feature into a random number (e.g. [f,g,b] becomes [1,2,3])

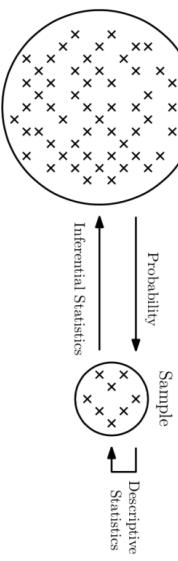
One-Hot Encoding: each of the m features becomes a vector of length m with containing only one 1 (e.g. [f, g, b] becomes [[1,0,0],[0,1,0],[0,0,1]])

Feature Hashing Scheme: turns arbitrary features into indices in a vector or matrix

Embeddings: if using words, convert words to vectors (word embeddings)

Statistical Analysis

Process of statistical reasoning: there is an underlying population of possible things we can potentially observe and only a small subset of them are actually sampled (ideally at random). Probability theory describes what properties our sample should have given the properties of the population, but *statistical inference* allows us to deduce what the full population is like after analyzing the sample.



Sampling From Distributions

Inverse Transform Sampling: Sampling points from a given probability distribution is sometimes necessary to run simulations or whether your data fits a particular distribution.

The general technique is called *inverse transform sampling* or Smirnov transform. First draw a random number p between [0,1]. Compute value x such that the CDF equals p : $F_X(x) = p$. Use x as the value to be the random value drawn from the distribution described by $F_X(x)$.

Monte Carlo Sampling: In higher dimensions, correctly sampling from a given distribution becomes more tricky. Generally want to use Monte Carlo methods, which typically follow these rules: define a domain of possible inputs, generate random inputs from a probability distribution over the domain, perform a deterministic calculation, and analyze the results.

Classic Statistical Distributions

Binomial Distribution (Discrete)

Assume X is distributed $\text{Bin}(n,p)$. X is the number of "successes" that we will achieve in n independent trials, where each trial is either a success or failure and each success occurs with the same probability p and each failure occurs with probability $q=1-p$.

$$\text{PDF: } P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$\text{EV: } \mu = np \quad \text{Variance} = npq$$

Normal/Gaussian Distribution (Continuous)

Assume X is distributed $\mathcal{N}(\mu, \sigma^2)$. It is a bell-shaped and symmetric distribution. Bulk of the values lie close to the mean and no value is too extreme. Generalization of the binomial distribution as $n \rightarrow \infty$.

$$\text{EV: } \mu \quad \text{Variance: } \sigma^2$$

$$\text{PDF: } P(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

Implications: 68%-95%-99% rule. 68% of probability mass fall within 1σ of the mean, 95% within 2σ , and 99.7% within 3σ .

Poisson Distribution (Discrete)

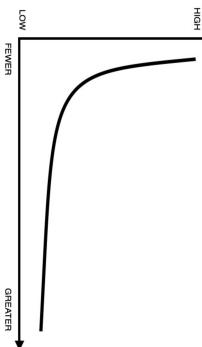
Assume X is distributed $\text{Pois}(\lambda)$. Poisson expresses the probability of a given number of events occurring in a fixed interval of time/space if these events occur independently and with a known constant rate λ .

$$\text{PDF: } P(x) = \frac{e^{-\lambda}\lambda^x}{x!} \quad \text{EV: } \lambda \quad \text{Variance} = \lambda$$

Power Law Distributions (Discrete)

Many data distributions have much longer tails than the normal or Poisson distributions. In other words, the change in one quantity varies as a *power* of another quantity. It helps measure the inequality in the world. e.g. wealth, word frequency and Pareto Principle (80/20 Rule)

$\text{PDF: } P(X=x) = cx^{-\alpha}$, where α is the law's exponent and c is the normalizing constant



Modeling- Overview

Modeling is the process of incorporating information into a tool which can forecast and make predictions. Usually, we are dealing with statistical modeling where we want to analyze relationships between variables. Formally, we want to estimate a function $f(X)$ such that:

$$Y = f(X) + \epsilon$$

where $X = (X_1, X_2, \dots, X_p)$ represents the input variables, Y represents the output variable, and ϵ represents random error.

Statistical learning is set of approaches for estimating this $f(X)$.

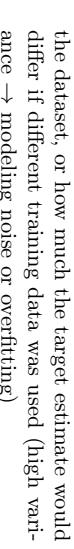
Why Estimate $f(X)$?

Prediction: once we have a good estimate $\hat{f}(X)$, we can use it to make predictions on new data. We treat \hat{f} as a black box, since we only care about the accuracy of the predictions, not why or how it works.

Inference: we want to understand the relationship between X and Y . We can no longer treat \hat{f} as a black box since we want to understand how Y changes with respect to $X = (X_1, X_2, \dots, X_p)$

More About ϵ

The error term ϵ is composed of the reducible and irreducible error, which will prevent us from ever obtaining a perfect \hat{f} estimate.



No Free Lunch Theorem No single machine learning algorithm is better than all the others on all problems. It is common to try multiple models and find one that works best for a particular problem.

Thinking Like Nate Silver

1. **Think Probabilistically** Probabilistic forecasts are more meaningful than concrete statements and should be reported as probability distributions (including σ) along with mean prediction μ .

2. **Incorporate New Information** Use live models, which continually updates using new information. To update, use Bayesian reasoning to calculate how probabilities change in response to new evidence.

3. **Look for Consensus Forecast** Use multiple distinct sources of evidence. Some models operate this way, such as boosting and bagging, which uses large number of weak classifiers to produce a strong one.

Modeling- Philosophies

Modeling is the process of incorporating information into a tool which can forecast and make predictions. Designing and validating models is important, as well as evaluating the performance of models. Note that the best forecasting model may not be the most accurate one.

Philosophies of Modeling

Occam's Razor Philosophical principle that the simplest explanation is the best explanation. In modeling, if we are given two models that predict equally well, we should choose the simpler one. Choosing the more complex one can often result in overfitting.

Bias Variance Trade-Off Inherent part of predictive modeling, where models with lower bias will have higher variance and vice versa. Goal is to achieve low bias and low variance.

- **Bias:** error from incorrect assumptions to make target function easier to learn (high bias \rightarrow missing relevant relations or underfitting)
- **Variance:** error from sensitivity to fluctuations in the dataset, or how much the target estimate would differ if different training data was used (high variance \rightarrow modeling noise or overfitting)

Modeling- Taxonomy

There are many different types of models. It is important to understand the trade-offs and when to use a certain type of model.

Parametric vs. Nonparametric

- **Parametric:** models that first make an assumption about a function form, or shape, of f (linear). Then fits the model. This reduces estimating f to just estimating set of parameters, but if our assumption was wrong, will lead to bad results.
- **Non-Parametric:** models that don't make any assumptions about f , which allows them to fit a wider range of shapes; but may lead to overfitting

Supervised vs. Unsupervised

- **Supervised:** models that fit input variables (x_1, x_2, \dots, x_n) to a known output variables $y_i = (y_1, y_2, \dots, y_n)$
- **Unsupervised:** models that take in input variables $x_i = (x_1, x_2, \dots, x_n)$, but they do not have an associated output to supervise the training. The goal is understand relationships between the variables or observations.

Blackbox vs. Descriptive

- **Blackbox:** models that make decisions, but we do not know what happens "under the hood" e.g. deep learning, neural networks
- **Descriptive:** models that provide insight into *why* they make their decisions e.g. linear regression, decision trees

First-Principle vs. Data-Driven

- **First-Principle:** models based on a prior belief of how the system under investigation works, incorporates domain knowledge (ad-hoc)
- **Data-Driven:** models based on observed correlations between input and output variables

Deterministic vs. Stochastic

- **Deterministic:** models that produce a single "prediction" e.g. yes or no, true or false
- **Stochastic:** models that produce probability distributions over possible events
- **Flat vs. Hierarchical**
 - Flat: models that solve problems on a single level, no notion of subproblems
 - Hierarchical: models that solve several different nested subproblems

Modeling- Evaluation Metrics

Need to determine how good our model is. Best way to assess models is out-of-sample predictions (data points your model has never seen).

Classification

| | Predicted Yes | Predicted No |
|------------|----------------------|----------------------|
| Actual Yes | True Positives (TP) | False Negatives (FN) |
| Actual No | False Positives (FP) | True Negatives (TN) |

Accuracy: ratio of correct predictions over total predictions. Misleading when class sizes are substantially different. $accuracy = \frac{TP+TN}{TP+TN+FN+FP}$

Precision: how often the classifier is correct when it predicts positive: $precision = \frac{TP}{TP+FP}$

Recall: how often the classifier is correct for all positive instances: $recall = \frac{TP}{TP+FN}$

F-Score: single measurement to describe performance: $F = 2 \cdot \frac{precision \cdot recall}{precision + recall}$

ROC Curves: plots true positive rates and false positive rates for various thresholds, or where the model determines if a data point is positive or negative (e.g. if >0.8 , classify as positive). Best possible area under the ROC curve (AUC) is 1, while random is 0.5, or the main diagonal line.

Regression

Errors are defined as the difference between a prediction \hat{y} and the actual result y .

Absolute Error: $\Delta = y' - y$

Squared Error: $\Delta^2 = (y' - y)^2$

Mean-Squared Error: $MSE = \frac{1}{n} \sum_{i=1}^n (y'_i - y_i)^2$

Root Mean-Squared Error: $RMSD = \sqrt{MSE}$

Absolute Error Distribution: Plot absolute error distribution: should be symmetric, centered around 0, bell-shaped, and contain rare extreme outliers.

Modeling- Evaluation Environment

Evaluation metrics provides use with the tools to estimate errors, but what should be the process to obtain the best estimate? Resampling involves repeatedly drawing samples from a training set and refitting a model to each sample, which provides us with additional information compared to fitting the model once, such as obtaining a better estimate for the test error.

Key Concepts

Training Data: data used to fit your models or the set used for learning

Validation Data: data used to tune the parameters of a model

Test Data: data used to evaluate how good your model is. Ideally your model should never touch this data until final testing/evaluation

Cross Validation

Class of methods that estimate test error by holding out a subset of training data from the fitting process.

Validation Set: split data into training set and validation set. Train model on training and estimate test error using validation. e.g. 80-20 split

Leave-One-Out CV (LOOCV): split data into training set and validation set, but the validation set consists of 1 observation. Then repeat n-1 times until all observations have been used as validation. Test error is the average of these n test error estimates.

k-Fold CV: randomly divide data into k groups (folds) of approximately equal size. First fold is used as validation and the rest as training. Then repeat k times and find average of the k estimates.

Bootstrapping: Methods that rely on random sampling with replacement. Bootstrapping helps with quantifying uncertainty associated with a given estimate or model.

Amplifying Small Data Sets

What can we do if we don't have enough data?

- **Create Negative Examples:** e.g. classifying pre-identical candidates, most people would be unqualified so label most as unqualified
- **Synthetic Data:** create additional data by adding noise to the real data

Linear Regression

Linear regression is a simple and useful tool for predicting a quantitative response. The relationship between input variables $\mathbf{X} = (X_1, X_2, \dots, X_p)$ and output variable Y takes the form:

$$Y \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

β_0, \dots, β_p are the unknown coefficients (parameters) which we are trying to determine. The best coefficients will lead us to the best "fit", which can be found by minimizing the *residual sum squares* (RSS), or the sum of the differences between the actual i th value and the predicted i th value. RSS = $\sum_{i=1}^n e_i$, where $e_i = y_i - \hat{y}_i$

How to find best fit?

Matrix Form: We can solve the closed-form equation for coefficient vector w : $w = (X^T X)^{-1} X^T Y$. X represents the input data and Y represents the output data. This method is used for smaller matrices, since inverting a matrix is computationally expensive.

Gradient Descent: First-order optimization algorithm.

We can find the minimum of a *convex* function by starting at an arbitrary point and repeatedly take steps in the downward direction, which can be found by taking the negative direction of the gradient. After several iterations, we will eventually converge to the minimum. In our case, the minimum corresponds to the coefficients with the minimum error, or the best line of fit. The learning rate α determines the size of the steps we take in the downward direction.

Gradient descent algorithm in two dimensions. Repeat until convergence.

1. $w_0^{t+1} := w_0^t - \alpha \frac{\partial}{\partial w_0} J(w_0, w_1)$
2. $w_1^{t+1} := w_1^t - \alpha \frac{\partial}{\partial w_1} J(w_0, w_1)$

For non-convex functions, gradient descent no longer guarantees an optimal solutions since there may be local minima. Instead, we should run the algorithm from different starting points and use the best local minima we find for the solution.

Stochastic Gradient Descent: instead of taking a step after sampling the *entire* training set, we take a small batch of training data at random to determine our next step. Computationally more efficient and may lead to faster convergence.

Linear Regression II

Improving Linear Regression

Subset/Feature Selection: approach involves identifying a subset of the p predictors that we believe to be best related to the response. Then we fit model using the reduced set of variables.

- Best, Forward, and Backward Subset Selection

Shrinkage/Regularization: all variables are used, but estimated coefficients are shrunk towards zero relative to the least squares estimate. λ represents the tuning parameter- as λ increases, flexibility decreases \rightarrow decreased variance but increased bias. The tuning parameter is key in determining the sweet spot between under and over-fitting. In addition, while Ridge will always produce a model with p variables, Lasso can force coefficients to be equal to zero.

- Lasso (L1): $\min \text{RSS} + \lambda \sum_{j=1}^p |\beta_j|$
- Ridge (L2): $\min \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$

Dimension Reduction: projecting p predictors into a M -dimensional subspace, where $M < p$. This is achieved by computing M different linear combinations of the variables. Can use PCA.

Miscellaneous: Removing outliers, feature scaling, removing multicollinearity (correlated variables)

Evaluating Model Accuracy

Residual Standard Error (RSE): $\text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}}$.

Generally, the smaller the better.

R^2 : Measure of fit that represents the proportion of variance explained, or the *variability in Y that can be explained using X* . It takes on a value between 0 and 1. Generally the higher the better. $R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$, where Total Sum of Squares (TSS) = $\sum (y_i - \bar{y})^2$

Evaluating Coefficient Estimates

Standard Error (SE) of the coefficients can be used to perform hypothesis tests on the coefficients:

H_0 : No relationship between X and Y , H_a : Some relationship exists. A p-value can be obtained and can be interpreted as follows: a small p-value indicates that a relationship between the predictor (X) and the response (Y) exists. Typical p-value cutoffs are around 5 or 1 %.

Logistic Regression

Logistic regression is used for classification, where the response variable is categorical rather than numerical.

The model works by predicting the probability that Y belongs to a particular category by first fitting the data to a linear regression model, which is then passed to the logistic function (below). The logistic function will always produce a S-shaped curve, so regardless of X , we can always obtain a sensible answer (between 0 and 1). If the probability is above a certain predetermined threshold (e.g. $P(\text{Yes}) > 0.5$), then the model will predict Yes.

$$p(X) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

How to find best coefficients?

Maximum Likelihood: The coefficients β_0, \dots, β_p are unknown and must be estimated from the training data. We seek estimates for β_0, \dots, β_p such that the predicted probability $\hat{p}(x_i)$ of each observation is a number close to one if its observed in a certain class and close to zero otherwise. This is done by maximizing the likelihood function:

$$L(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i':y_{i'}=1} (1 - p(x_i))$$

Potential Issues

Imbalanced Classes: imbalance in classes in training data lead to poor classifiers. It can result in a lot of false positives and also lead to few training data. Solutions include forcing balanced data by removing observations from the larger class, replicate data from the smaller class, or heavily weigh the training examples toward instances of the larger class.

Multi-Class Classification: the more classes you try to predict, the harder it will be for the classifier to be effective. It is possible with logistic regression, but another approach, such as Linear Discriminant Analysis (LDA), may prove better.

Distance/Network Methods

Interpreting examples as points in space provides a way to find natural groupings or clusters among data e.g. which stars are the closest to our sun? Networks can also be built from point sets (vertices) by connecting related points.

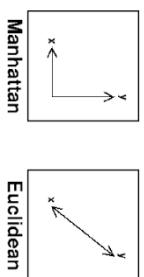
Measuring Distances/Similarity Measure

There are several ways of measuring distances between points a and b in d dimensions- with closer distances implying similarity.

Minkowski Distance Metric: $d_k(a, b) = \sqrt[k]{\sum_{i=1}^d |a_i - b_i|^k}$

The parameter k provides a way to tradeoff between the largest and the total dimensional difference. In other words, larger values of k place more emphasis on large differences between feature values than smaller values. Selecting the right k can significantly impact the the meaningfulness of your distance function. The most popular values are 1 and 2.

- Manhattan ($k=1$): city block distance, or the sum of the absolute difference between two points
- Euclidean ($k=2$): straight line distance



Weighted Minkowski: $d_k(a, b) = \sqrt[k]{\sum_{i=1}^d w_i |a_i - b_i|^k}$, in some scenarios, not all dimensions are equal. Can convey this idea using w_i . Generally not a good idea- should normalize data by Z-scores before computing distances.

Cosine Similarity: $\cos(a, b) = \frac{a \cdot b}{\|a\| \|b\|}$, calculates the similarity between 2 non-zero vectors, where $a \cdot b$ is the dot product (normalized between 0 and 1); higher values imply more similar vectors

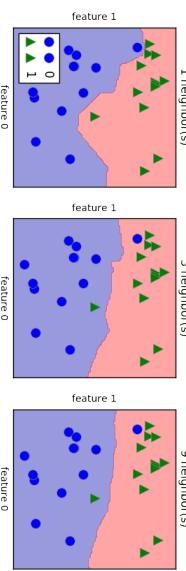
Kullback-Leibler Divergence: $KL(A||B) = \sum_{i=1}^d a_i \log_2 \frac{a_i}{b_i}$ KL divergence measures the distances between probability distributions by measuring the uncertainty gained or uncertainty lost when replacing a distribution A with distribution B. However, this is not a metric but forms the basis for the Jensen-Shannon Divergence Metric. **Jensen-Shannon:** $JS(A, B) = \frac{1}{2}KL(A||M) + \frac{1}{2}KL(M||B)$, where M is the average of A and B . The JS function is the right metric for calculating distances between probability distributions

Nearest Neighbor Classification

Distance functions allow us to identify the points closest to a given target, or the *nearest neighbors (NN)* to a given point. The advantages of NN include simplicity, interpretability and non-linearity.

k-Nearest Neighbors

Given a positive integer k and a point x_0 , the KNN classifier first identifies k points in the training data most similar to x_0 , then estimates the conditional probability of x_0 being in class j as the fraction of the k points whose values belong to j . The optimal value for k can be found using cross validation.



KNN Algorithm

1. Compute distance $D(a, b)$ from point b to all points
2. Select k closest points and their labels
3. Output class with most frequent labels in k points

Optimizing KNN

Comparing a query point a in d dimensions against n training examples computes with a runtime of $O(nd)$, which can cause lag as points reach millions or billions. Popular choices to speed up KNN include:

- **Voronoi Diagrams:** partitioning plane into regions based on distance to points in a specific subset of the plane
- **Grid Indexes:** carve up space into d -dimensional boxes or grids and calculate the NN in the same cell as the point
- **Locality Sensitive Hashing (LSH):** abandons the idea of finding the exact nearest neighbors. Instead, batch up nearby points to quickly find the most appropriate bucket B for our query point. LSH is defined by a hash function $h(p)$ that takes a point/vector as input and produces a number/ code as output, such that it is likely that $h(a) = h(b)$ if a and b are close to each other, and $h(a) \neq h(b)$ if they are far apart.

Clustering

Clustering is the problem of grouping points by similarity using distance metrics, which ideally reflect the similarities you are looking for. Often items come from logical "sources" and clustering is a good way to reveal those origins. Perhaps the first thing to do with any data set. Possible applications include: hypothesis development, modeling over smaller subsets of data, data reduction, outlier detection.

K-Means Clustering

Simple and elegant algorithm to partition a dataset into K distinct, non-overlapping clusters.

1. Choose a K . Randomly assign a number between 1 and K to each observation. These serve as initial cluster assignments
2. Iterate until cluster assignments stop changing
 - (a) For each of the K clusters, compute the cluster centroid. The k th cluster centroid is the vector of the p feature means for the observations in the k th cluster.
 - (b) Assign each observation to the cluster whose centroid is closest (where closest is defined using distance metric).

Since the results of the algorithm depends on the initial random assignments, it is a good idea to repeat the algorithm from different random initializations to obtain the best overall results. Can use MSE to determine which cluster assignment is better.

Hierarchical Clustering

Alternative clustering algorithm that does not require us to commit to a particular K . Another advantage is that it results in a nice visualization called a **dendrogram**. Observations that fuse at bottom are similar, where those at the top are quite different- we draw conclusions based on the location on the vertical rather than horizontal axis.

1. Begin with n observations and a measure of all the $\frac{(n)(n-1)}{2}$ pairwise dissimilarities. Treat each observation as its own cluster.
2. For $i = n, n-1, \dots, 2$
 - (a) Examine all pairwise inter-cluster dissimilarities among the i clusters and identify the pair of clusters that are least dissimilar (most similar). Fuse these two clusters. The dissimilarity between these two clusters indicates height in dendrogram where fusion should be placed.
 - (b) Assign each observation to the cluster whose centroid is closest (where closest is defined using distance metric).

Linkage: Complete (max dissimilarity), Single (min), Average, Centroid (between centroids of cluster A and B)

Machine Learning Part I

Comparing ML Algorithms

Power and Expressibility: ML methods differ in terms of complexity. Linear regression fits linear functions while NN define piecewise-linear separation boundaries. More complex models can provide more accurate models, but at the risk of overfitting.

Interpretability: some models are more transparent and understandable than others (white box vs. black box models)

Ease of Use: some models feature few parameters/decisions (linear regression/NN), while others require more decision making to optimize (SVMs)

Prediction Speed: models differ in how fast they fit the necessary parameters predictions given a query

| Method | Power of Expression | Ease of Interpretation | Ease of Use | Training Speed | Prediction Speed |
|-------------------------|---------------------|------------------------|-------------|----------------|------------------|
| Linear Regression | 5 | 9 | 9 | 9 | — |
| Nearest Neighbor | 5 | 9 | 8 | 10 | 2 |
| Naive Bayes | 4 | 8 | 7 | 9 | 8 |
| Decision Trees | 8 | 8 | 7 | 7 | 9 |
| Support Vector Machines | 9 | 6 | 6 | 7 | 7 |
| Boosting | 9 | 6 | 6 | 6 | 6 |
| Graphical Models | 9 | 8 | 3 | 4 | 4 |
| Deep Learning | 10 | 3 | 4 | 3 | 7 |

Naive Bayes

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the "naive" assumption of independence between every pair of features.

Problem: Suppose we need to classify vector $X = x_1 \dots x_n$ into m classes, $C_1 \dots C_m$. We need to compute the probability of each possible class given X , so we can assign X the label of the class with highest probability. We can calculate a probability using the Bayes' Theorem:

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)}$$

Where:

1. $P(C_i)$: the prior probability of belonging to class i
2. $P(X)$: normalizing constant, or probability of seeing the given input vector over all possible input vectors
3. $P(X|C_i)$: the conditional probability of seeing input vector X given we know the class is C_i

The prediction model will formally look like:

$$C(X) = \arg\max_{i \in \text{classes}} \frac{P(X|C_i)P(C_i)}{P(X)}$$

where $C(X)$ is the prediction returned for input X .

Machine Learning Part II

Decision Trees

Binary branching structure used to classify an arbitrary input vector X . Each node in the tree contains a simple feature comparison against some field ($x_i > 42?$). Result of each comparison is either true or false, which determines if we should proceed along to the left or right child of the given node. Also known as sometimes called classification and regression trees (CART).

Person



Advantages: Non-linearity, support for categorical variables, easy to interpret, application to regression.

Disadvantages: Prone to overfitting, unstable (not robust to noise), high variance, low bias

Note: rarely do models just use one decision tree. Instead, we aggregate many decision trees using methods like ensembling, bagging, and boosting

Ensembles, Bagging, Random Forests, Boosting

Ensemble learning is the strategy of combining many different classifiers/models into one predictive model. It revolves around the idea of voting: a so-called "wisdom of crowds" approach. The most predicted class will be the final prediction.

Bagging: ensemble method that works by taking B bootstrap samples of the training data and constructing B trees, each tree training on a distinct subsample as

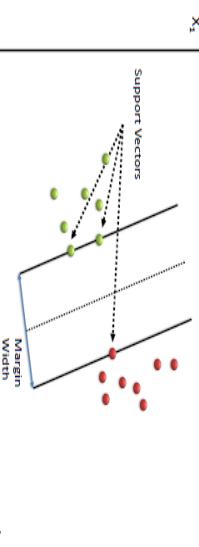
Random Forests: builds on bagging by decorrelating the trees. We do everything the same like in bagging, but when we build the trees, everytime we consider a split, a random sample of the p predictors is chosen as split candidates, not the full set (typically $m \approx \sqrt{p}$). When $m = p$, then we are just doing bagging.

Boosting: the main idea is to improve our model where it is not performing well by using information from previously constructed classifiers. Slow learner. Has 3 tuning parameters: number of classifiers B , learning parameter λ , interaction depth d (controls interaction order of model).

Machine Learning Part III

Support Vector Machines

Work by constructing a hyperplane that separates points between two classes. The hyperplane is determined using the maximal margin hyperplane, which is the hyperplane that is the maximum distance from the training observations. This distance is called the margin. Points that fall on one side of the margin are classified as -1 and the other +1.

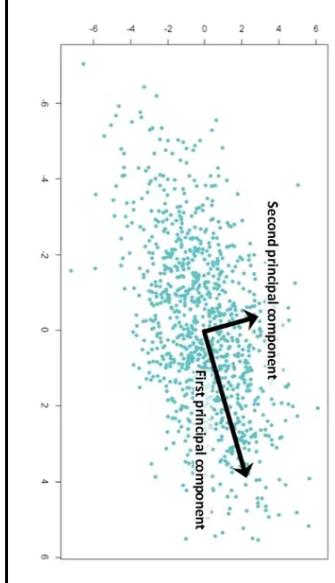


Principal Component Analysis (PCA)

Principal components allow us to summarize a set of correlated variables with a smaller set of variables that collectively explain most of the variability in the original set. Essentially, we are "dropping" the least important feature variables.

Principal Component Analysis

which principal components are calculated and the use of them to analyzing and understanding the data. PCA is an unsupervised approach and is used for dimensionality reduction, feature extraction, and data visualization. Variables after performing PCA are independent. Scaling variables is also important while performing PCA.



Machine Learning Part IV

ML Terminology and Concepts

Deep Learning Part I

What is Deep Learning?

Features: input data/variables used by the ML model
Feature Engineering: transforming input features to be more useful for the models, e.g. mapping categories to buckets, normalizing between -1 and 1, removing null

Train/Eval/Test: training is data used to optimize the model, evaluation is used to assess the model on new data during training, test is used to provide the final result

Classification/Regression: regression is prediction a number (e.g. housing price), classification is prediction from a set of categories(e.g. predicting red/blue/green)

Linear Regression: predicts an output by multiplying and summing input features with weights and biases

Logistic Regression: similar to linear regression but predicts a probability

Oversfitting: model performs great on the input data but poorly on the test data (combat by dropout, early stopping, or reduce # of nodes or layers)

Bias/Variance: how much output is determined by the features, more variance often can mean overfitting, more bias can mean a bad model

Regularization: variety of approaches to reduce overfitting, including adding the weights to the loss function, randomly dropping layers (dropout)

Ensemble Learning: training multiple models with different parameters to solve the same problem

A/B testing: statistical way of comparing 2+ techniques to determine which technique performs better and also if difference is statistically significant

Baseline Model: simple model/heuristic used as reference point for comparing how well a model is performing

Bias: prejudice or favoritism towards some things, people, or groups over others that can affect collection/sampling and interpretation of data, the design of a system, and how users interact with a system

Dynamic Model: model that is trained online in a continuously updating fashion

Static Model: model that is trained offline

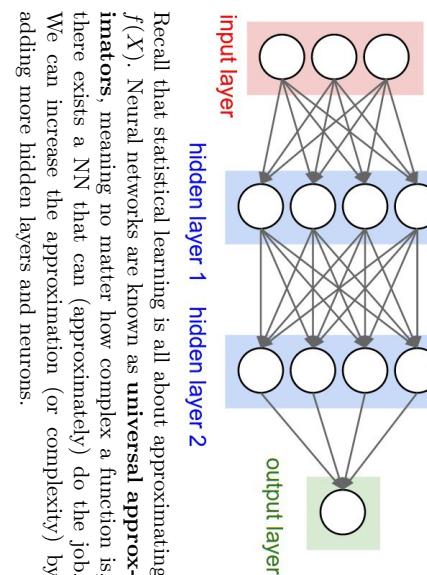
Normalization: process of converting an actual range of values into a standard range of values, typically -1 to +1

Independently and Identically Distributed (i.i.d): data drawn from a distribution that doesn't change, and where each value drawn doesn't depend on previously drawn values; ideal but rarely found in real life

Hyperparameters: the "knobs" that you tweak during successive runs of training a model

Generalization: refers to a model's ability to make correct predictions on new, previously unseen data as opposed to the data used to train the model

Cross-Entropy: quantifies the difference between two probability distributions



Recall that statistical learning is all about approximating $f(X)$. Neural networks are known as **universal approximators**, meaning no matter how complex a function is, there exists a NN that can (approximately) do the job. We can increase the approximation (or complexity) by adding more hidden layers and neurons.

Popular Architectures

There are different kinds of NNs that are suitable for certain problems, which depend on the NN's architecture.

Linear Classifier: takes input features and combines them with weights and biases to predict output value

DNN: deep neural net, contains intermediate layers of nodes that represent "hidden features" and activation functions to represent non-linearity

CNN: convolutional NN, has a combination of convolutional, pooling, dense layers, popular for image classification.

Transfer Learning: use existing trained models as starting points and add additional layers for the specific use case. idea is that highly trained existing models know general features that serve as a good starting point for training a small network on specific examples

RNN: recurrent NN, designed for handling a sequence of inputs that have "memory" of the sequence. LSTMs are a fancy version of RNNs, popular for NLP

GAN: general adversarial NN, one model creates fake examples, and another model is served both fake example and real examples and is asked to distinguish

Wide and Deep: combines linear classifiers with deep neural net classifiers, "wide" linear parts represent memorizing specific examples and "deep" parts represent understanding high level features

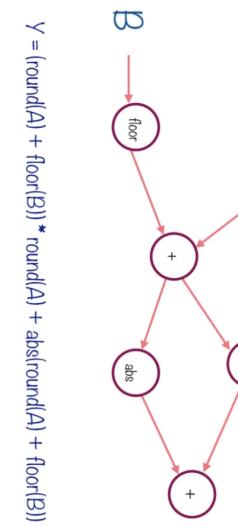
Deep Learning Part II

Tensorflow

Tensorflow is an open source software library for numerical computation using data flow graphs. Everything in TF is a graph, where nodes represent operations on data and edges represent the data. Phase 1 of TF is building up a computation graph and phase 2 is executing it. It is also distributed, meaning it can run on either a cluster of machines or just a single machine.

TF is extremely popular/suitable for working with Neural Networks, since the way TF sets up the computational graph pretty much resembles a NN.

Tensors Flow Through the Graph



Tensors

In a graph, tensors are the edges and are multidimensional data arrays that flow through the graph. Central unit of data in TF and consists of a set of primitive values shaped into an array of any number of dimensions. A tensor is characterized by its rank (# dimensions in tensor), shape (# of dimensions and size of each dimension), data type (data type of each element in tensor).

Placeholders and Variables

Variables: best way to represent shared, persistent state manipulated by your program. These are the parameters of the ML model are altered/trained during the training process. Training variables.

Placeholders: way to specify inputs into a graph that hold the place for a Tensor that will be fed at runtime. They are assigned once, do not change after. Input nodes

Deep Learning Part III

Deep Learning Terminology and Concepts

Neuron: node in a NN, typically taking in multiple input values and generating one output value, calculates the output value by applying an activation function (nonlinear transformation) to a weighted sum of input values

Weights: edges in a NN, the goal of training is to determine the optimal weight for each feature; if weight = 0, corresponding feature does not contribute

Neural Network: composed of neurons (simple building blocks that actually “learn”), contains activation functions that makes it possible to predict non-linear outputs

Activation Functions: mathematical functions that introduce non-linearity to a network e.g. RELU, tanh

Sigmoid Function: function that maps very negative numbers to a number very close to 0, huge numbers close to 1, and 0 to .5. Useful for predicting probabilities

Gradient Descent/Backpropagation: fundamental loss optimizer algorithms, of which the other optimizers are usually based. Backpropagation is similar to gradient descent but for neural nets

Optimizer: operation that changes the weights and biases to reduce loss e.g. Adagrad or Adam

Weights / Biases: weights are values that the input features are multiplied by to predict an output value. Biases are the value of the output given a weight of 0.

Converge: algorithm that converges will eventually reach an optimal answer, even if very slowly. An algorithm that doesn't converge may never reach an optimal answer.

Learning Rate: rate at which optimizers change weights and biases. High learning rate generally trains faster but risks not converging, whereas a lower rate trains slower

Numerical Instability: issues with very large/small values due to limits of floating point numbers in computers

Embeddings: mapping from discrete objects, such as words, to vectors of real numbers, useful because classifiers/neural networks work well on vectors of real numbers

Convolutional Layer: series of convolutional operations, each acting on a different slice of the input matrix. **Dropout:** method for regularization in training NNs, works by removing a random selection of some units in a network layer for a single gradient step

Early Stopping: method for regularization that involves ending model training early

Gradient Descent: technique to minimize loss by computing the gradients of loss with respect to the model's parameters, conditioned on training data

Pooling: Reducing a matrix (or matrices) created by an earlier convolutional layer to a smaller matrix. Pooling usually involves taking either the maximum or average value across the pooled area

Big Data- Hadoop Overview

Data can no longer fit in memory on one machine (monolithic), so a new way of computing was devised using a group of computers to process this “big data” (distributed). Such a group is called a cluster, which makes up server farms. All of these servers have to be coordinated in the following ways: partition data, coordinate computing tasks, handle fault tolerance/recovery, and allocate capacity to process.

Hadoop

Hadoop is an open source *distributed* processing framework that manages data processing and storage for big data applications running in clustered systems. It is comprised of 3 main components:

- **Hadoop Distributed File System (HDFS):** a distributed file system that provides high-throughput access to application data by partitioning data across many machines
- **YARN:** framework for job scheduling and cluster resource management (task coordination)
- **MapReduce:** YARN-based system for parallel processing of large data sets on multiple machines

HDFS

Each disk on a different machine in a cluster is comprised of 1 master node and the rest are workers/data nodes;

The **master node** manages the overall file system by storing the directory structure and the metadata of the files. The **data nodes** physically store the data. Large files are broken up and distributed across multiple machines, which are also replicated across multiple machines to provide fault tolerance.

MapReduce

Parallel programming paradigm which allows for processing of huge amounts of data by running processes on multiple machines. Defining a MapReduce job requires two stages: map and reduce.

- **Map:** operation to be performed in parallel on small portions of the dataset. the output is a key-value pair $< K, V >$
- **Reduce:** operation to combine the results of Map

YARN- Yet Another Resource Negotiator

Coordinates tasks running on the cluster and assigns new nodes in case of failure. Comprised of 2 subcomponents: the resource manager and the node manager. The **resource manager** runs on a single master node and schedules tasks across nodes. The **node manager** runs on all other nodes and manages tasks on the individual node.

Big Data- Hadoop Ecosystem

An entire ecosystem of tools have emerged around Hadoop, which are based on interacting with HDFS. Below are some popular ones:

Hive: data warehouse software built on top of Hadoop that facilitates reading, writing, and managing large datasets residing in distributed storage using SQL-like queries (HiveQL). Hive abstracts away underlying MapReduce jobs and returns HDFS in the form of tables (not HDFS).

Pig: high level scripting language (Pig Latin) that enables writing complex data transformations. It pulls unstructured/incomplete data from sources, cleans it, and places it in a database/data warehouses. Pig performs ETL into data warehouse while Hive queries from data warehouse to perform analysis (GCP: DataFlow).

Spark: framework for writing fast, distributed programs for data processing and analysis. Spark solves similar problems as Hadoop Map Reduce but with a fast in-memory approach. It is an unified engine that supports SQL queries, streaming data, machine learning and graph processing. Can operate separately from Hadoop but integrates well with Hadoop. Data is processed using Resilient Distributed Datasets (RDDs), which are immutable, lazily evaluated, and tracks lineage.

Hbase: non-relational, NoSQL, column-oriented database management system that runs on top of HDFS. Well suited for sparse data sets (GCP: BigTable)

Flink/Kafka: stream processing framework. Batch streaming is for bounded, finite datasets, with periodic updates, and delayed processing. Stream processing is for unbounded datasets, with continuous updates, and immediate processing. Stream data and stream processing must be decoupled via a message queue. Can group streaming data (windows) using tumbling (non-overlapping time), sliding (overlapping time), or session (session gap) windows.

Beam: programming model to define and execute data processing pipelines, including ETL, batch and stream (continuous) processing. After building the pipeline, it is executed by one of Beam's distributed processing back-ends (Apache Apex, Apache Flink, Apache Spark, and Google Cloud Dataflow). Modeled as a Directed Acyclic Graph (DAG).

Oozie: workflow scheduler system to manage Hadoop jobs

Sqoop: transferring framework to transfer large amounts of data into HDFS from relational databases (MySQL) source manager runs on a single master node and schedules tasks across nodes. The node manager runs on all other nodes and manages tasks on the individual node.

SQL Part I

Structured Query Language (SQL) is a declarative language used to access & manipulate data in databases. Usually the database is a Relational Database Management System (RDBMS), which stores data arranged in relational database tables. A table is arranged in columns and rows, where columns represent characteristics of stored data and rows represent actual data entries.

Basic Queries

- filter columns: **SELECT** col1, col3... **FROM** table1
- filter the rows: **WHERE** col4 = 1 **AND** col5 = 2
- aggregate the data: **GROUP BY**...
- limit aggregated data: **HAVING count(*) > 1**
- order of the results: **ORDER BY** col2

Useful Keywords for **SELECT**

DISTINCT- return unique results
BETWEEN a AND b- limit the range, the values can be numbers, text, or dates

LIKE- pattern search within the column text
IN (a, b, c) - check if the value is contained among given

Data Modification

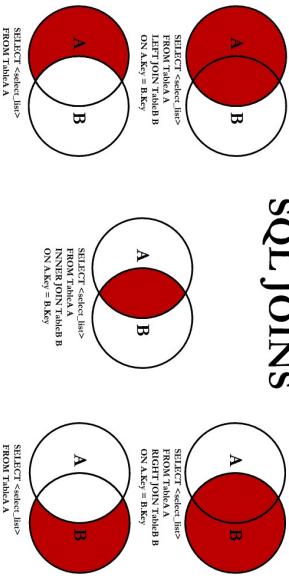
- update specific data with the **WHERE** clause:
- UPDATE** table1 **SET** col1 = 1 **WHERE** col2 = 2
- insert values manually

INSERT INTO table1 (col1,col3) **VALUES** (val1,val3);
- by using the results of a query
INSERT INTO table1 (col1,col3) **SELECT** col1,col2
FROM table2;

Joins

The JOIN clause is used to combine rows from two or more tables, based on a related column between them.

SQL JOINS



```
SELECT <select list>
  FROM TableA
    LEFT JOIN TableB
      ON A.Key = B.Key
    FULL OUTER JOIN TableB
      ON A.Key = B.Key
 WHERE A.Key IS NULL;
```

© CL Mönch 2008

Python- Data Structures

Data structures are a way of storing and manipulating data and each data structure has its own strengths and weaknesses. Combined with algorithms, data structures allow us to efficiently solve problems. It is important to know the main types of data structures that you will need to efficiently solve problems.

Lists: or arrays, ordered sequences of objects, mutable

```
>>> l = [42, 3.14, "hello", "world"]
```

Tuples: like lists, but immutable

```
>>> t = (42, 3.14, "hello", "world")
```

Dictionaries: hash tables, key-value pairs, unsorted queues; supports append, appendleft, pop, rotate, etc

```
>>> s = deque([42, 3.14, "hello", "world"])
```

Sets: mutable, unordered sequence of unique elements. frozensets are just immutable sets

```
>>> s = set([42, 3.14, "hello", "world"])
```

Collections Module

deque: double-ended queue, generalization of stacks and queues; supports append, appendleft, pop, rotate, etc

```
>>> s = deque([42, 3.14, "hello", "world"])
```

Counter: dict subclass, unordered collection where elements are stored as keys and counts stored as values

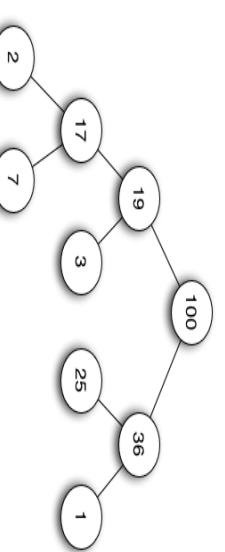
```
>>> c = Counter('apple')
>>> print(c)
Counter({'p': 2, 'a': 1, 'l': 1, 'e': 1})
```

heapq Module

Heap Queue: priority queue, heaps are binary trees for which every parent node has a value greater than or equal to any of its children (max-heap), order is important; supports push, pop, pushpop, heapify, replace functionality

```
>>> heap = []
>>> for n in data:
...     heappush(heap, n)
...     ...
>>> heap
```

```
[0, 1, 3, 6, 2, 8, 4, 7, 9, 5]
```



DataCamp

Data Science Cheat Sheet for Business Leaders

Data Science Basics

Types of Data Science

- **Descriptive Analytics (Business Intelligence):** Get useful data in front of the right people in the form of dashboards, reports, and emails
 - Which customers have churned?
 - Which homes have sold in a given location, and do homes of a certain size sell more quickly?
- **Predictive Analytics (Machine Learning):** Put data science models continuously into production
 - Which customers may churn?
 - How much will a home sell for, given its location and number of rooms?

- **Prescriptive Analytics (Decision Science):** Use data to help a company make decisions

- What should we do about the particular types of customers that are prone to churn?
- How should we market a home to sell quickly, given its location and number of rooms?

The Standard Data Science Workflow

- 1 **Data Collection:** Compile data from different sources and store it for efficient access



- 2 **Exploration and Visualization:** Explore and visualize data through dashboards



- 3 **Experimentation and Prediction:** The buzziest topic in data science—machine learning!



datacamp.com/courses/data-science-for-business-leaders



datacamp.com/business

Building a Data Science Team

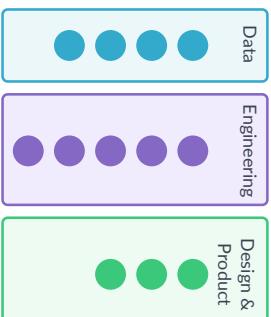
Your data team members require different skills for different purposes.

| Data Engineer | Data Analyst | Machine Learning Engineer | Data Scientist |
|-------------------------|-------------------------------|--|---|
| Store and maintain data | Visualize and describe data | Write production-level code to predict with data | Build custom models to drive business decisions |
| SQL/Java/Scala/Python | SQL + BI Tools + Spreadsheets | Python/Java/R | Python/R/SQL |

Data Science Team Organizational Models

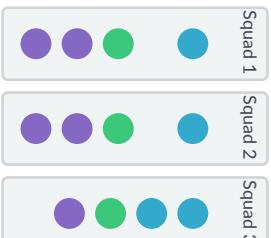
Centralized/Isolated

The data team is the owner of data and answers requests from other teams



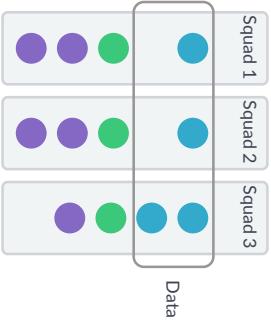
Embedded

Data experts are dispersed across an organization and report to functional leaders



Hybrid

Data experts sit with functional teams and also report to the Chief Data Scientist—so data is an organizational priority



Python 3 Cheat Sheet

Latest version on :
<https://perso.limsi.fr/pointal/python:memento>

| Base Types | | Container Types | |
|--|---|--|-------------------------|
| integer, float, boolean, string, bytes | | ordered sequences, fast index access, repeatable values | |
| int 783 0 -192 0b010 0o642 0xF3 | zero binary octal hexa | list [1, 5, 9] ["x", 11, 8.9] ["mot"] | [] |
| float 9.23 0.0 -1.7e-6 | | tuple (1, 5, 9) (11, "y", 7.4) ("mot",) | () |
| bool True False | x10 ⁻⁶ | Non modifiable values (immutables) | "::" |
| str "One\nTwo" | Multiline string: '''X\y\tz''' 1\t2\t3''' | expression with only commas → tuple | b''' |
| escaped new line | escaped tab | str bytes (ordered sequences of chars / bytes) | |
| 'I\'m' | | | |
| escaped ' | | | |
| bytes b'toto\xfe\x77' | hexadecimal octal | key containers, no a priori order, fast key access, each key is unique | {} |
| | immutable | dictionary dict {"key": "value"} dict(a=3, b=4, k="v") | {} |
| | | (key/value associations) {1: "one", 3: "three", 2: "two", 3.14: "pi"} | |
| | | collection set {"key1", "key2"} {1, 9, 3, 0} | {} |
| | | keys=hashable values (base types, immutables...) | frozenset immutable set |
| | | | empty |

| Identifiers |
|---|
| for variables, functions, modules, classes... names |
| a...zA...Z_ followed by a...zA...Z_0...9 |
| ▫ diacritics allowed but should be avoided |
| ▫ language keywords forbidden |
| ▫ lower/UPPER case discrimination |
| ◎ a toto x7 y_max BigOne |
| ◎ by and for |

| Variables assignment |
|--|
| ▫ assignment ⇔ binding of a name with a value |
| 1) evaluation of right side expression value |
| 2) assignment in order with left side names |
| x=1.2+8+sin(y) |
| a=b=c=0 assignment to same value |
| y, z, r=9.2, -7.6, 0 multiple assignments |
| a, b=b, a values swap |
| a, *b=seq unpacking of sequence in *a, b=seq item and list |
| x+=3 increment ⇔ x=x+3 |
| x-=2 decrement ⇔ x=x-2 |
| x=None « undefined » constant value |
| del x remove name x |
| and |
| * |
| /= |
| %= |
| ... and |

| Conversions |
|--|
| int("15") → 15 |
| int("3f", 16) → 63 |
| int(15.56) → 15 |
| float("-11.24e8") → -1124000000.0 |
| round(15.56, 1) → 15.6 |
| bool(x) False for null x, empty container x, None or False x; True for other x |
| str(x) → "..." representation string of x for display (cf. formatting on the back) |
| chr(64) → '@' ord('@') → 64 code ↔ char |
| repr(x) → "..." literal representation string of x |
| bytes([72, 9, 64]) → b'H\t@' |
| list("abc") → ['a', 'b', 'c'] |
| dict([(3, "three"), (1, "one")]) → {1: 'one', 3: 'three'} |
| set(["one", "two"]) → {'one', 'two'} |
| separator str and sequence of str → assembled str |
| ':'.join(['toto', '12', 'pswd']) → 'toto:12:pswd' |
| str splitted on whitespaces → list of str |
| "words with spaces".split() → ['words', 'with', 'spaces'] |
| str splitted on separator str → list of str |
| "1,4,8,2".split(",") → ['1', '4', '8', '2'] |
| sequence of one type → list of another type (via list comprehension) |
| [int(x) for x in ('1', '29', '-3')] → [1, 29, -3] |

| for lists, tuples, strings, bytes... | Sequence Containers Indexing | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|--|--|----|----|----|----|----|----------------|---|---|---|---|---|--------------------------|----|----|----|----|----|----------------|---|---|---|---|---|----------------|----|----|----|----|----|---|
| <table border="1"> <tr> <td>negative index</td><td>-5</td><td>-4</td><td>-3</td><td>-2</td><td>-1</td></tr> <tr> <td>positive index</td><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td></tr> <tr> <td>lst=[10, 20, 30, 40, 50]</td><td>10</td><td>20</td><td>30</td><td>40</td><td>50</td></tr> <tr> <td>positive slice</td><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td></tr> <tr> <td>negative slice</td><td>-5</td><td>-4</td><td>-3</td><td>-2</td><td>-1</td></tr> </table> | negative index | -5 | -4 | -3 | -2 | -1 | positive index | 0 | 1 | 2 | 3 | 4 | lst=[10, 20, 30, 40, 50] | 10 | 20 | 30 | 40 | 50 | positive slice | 0 | 1 | 2 | 3 | 4 | negative slice | -5 | -4 | -3 | -2 | -1 | Individual access to items via lst[index] |
| negative index | -5 | -4 | -3 | -2 | -1 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| positive index | 0 | 1 | 2 | 3 | 4 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| lst=[10, 20, 30, 40, 50] | 10 | 20 | 30 | 40 | 50 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| positive slice | 0 | 1 | 2 | 3 | 4 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| negative slice | -5 | -4 | -3 | -2 | -1 | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Access to sub-sequences via lst[start slice:end slice:step] | lst[0] → 10 → first one lst[1] → 20 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| lst[:1] → [10, 20, 30, 40] lst[::1] → [50, 40, 30, 20, 10] lst[1:-1] → [20, 30, 40] lst[::2] → [10, 30, 50] | lst[-1] → 50 → last one lst[-2] → 40 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| lst[1:-1] → [20, 30, 40] lst[::2] → [10, 30, 50] | On mutable sequences (list), remove with del lst[3] and modify with assignment lst[4]=25 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| lst[1:-1] → [20, 30, 40] lst[::2] → [10, 30, 50] | shallow copy of sequence | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Missing slice indication → from start / up to end. | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| On mutable sequences (list), remove with del lst[3:5] and modify with assignment lst[1:4]=[15, 25] | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

| Boolean Logic | Statements Blocks | Modules/Names Imports |
|--|--------------------------------------|---|
| Comparisons : < > <= >= == != (boolean results) ≤ ≥ = ≠ | Items count | module truc⇒file truc.py |
| a and b logical and both simultaneously | len(lst) → 5 | from monmod import nom1, nom2 as fct |
| a or b logical or one or other or both | ▫ index from 0 (here from 0 to 4) | → direct access to names, renaming with as |
| ▫ pitfall : and and or return value of a or of b (under shortcut evaluation). ⇒ ensure that a and b are booleans. | | import monmod → access via monmod.nom1 ... |
| not a logical not | | ▫ modules and packages searched in python path (cf. sys.path) |
| True False | True and False constants | |

| Floating numbers... approximated values | Maths | Conditional Statement |
|---|---|---|
| Operators: + - * / // % ** Priority (...) × ÷ ↑ ↑ a ^b Priority (...) integer ÷ remainder @ → matrix × python3.5+numpy (1+5.3)*2→12.6 abs(-3.2)→3.2 round(3.57,1)→3.6 pow(4,3)→64.0 ▫ usual order of operations | angles in radians from math import sin, pi... sin(pi/4)→0.707... cos(2*pi/3)→-0.4999... sqrt(81)→9.0 log(e**2)→2.0 ceil(12.5)→13 floor(12.5)→12 modules math, statistics, random, decimal, fractions, numpy, etc. (cf. doc) | statement block executed only if a condition is true |
| | | if logical condition: → statements block |
| | | can go with several elif, elif... and only one final else. Only the block of first true condition is executed. |
| | | ▫ with a var x: if bool(x)==True: → if x: if bool(x)==False: → if not x: |
| | | Signalizing an error: raise ExcClass(...) |
| | | Errors processing: try: → normal processing block |
| | | except Exception as e: → error processing block |
| | | Exceptions on Errors ▫ normal raise X(processing error processing error raise processing ▫ finally block for final processing in all cases. |

Conditional Loop Statement

statements block executed as long as condition is true

while logical condition:

→ statements block

Loop Control

break immediate exit
continue next iteration
else block for normal loop exit.

Algo: $i=100$
 $S = \sum_{i=1}^{100} i^2$

Iterative Loop Statement

statements block executed for each item of a container or iterator

for var in sequence:

→ statements block

next → finish

Display

print("v=", 3, "cm : ", x, ", ", y+4)

items to display : literal values, variables, expressions

Input

print options:

- sep=" " items separator, default space
- end="\n" end of print, default new line
- file=sys.stdout print to file, default standard output

s = input("Instructions:")

input always returns a string, convert it to required type (cf. boxed Conversions on the other side).

Generic Operations on Containers

len(c) → items count
min(c) max(c) sum(c) Note: For dictionaries and sets, these operations use keys.
sorted(c) → list sorted copy
val in c → boolean, membership operator in (absence not in)
enumerate(c) → iterator on (index, value)
zip(c1, c2...) → iterator on tuples containing c1 items at same index
all(c) → True if all c items evaluated to true, else False
any(c) → True if at least one item of c evaluated true, else False
Specific to ordered sequences containers (lists, tuples, strings, bytes...)
reversed(c) → inverted iterator c*5 → duplicate c+c2 → concatenate
c.index(val) → position c.count(val) → events count
import copy
copy.copy(c) → shallow copy of container
copy.deepcopy(c) → deep copy of container

Operations on Lists

modify original list

lst.append(val) add item at end
lst.extend(seq) add sequence of items at end
lst.insert(idx, val) insert item at index
lst.remove(val) remove first item with value val
lst.pop(idx) → value remove & return item at index idx (default last)
lst.sort() lst.reverse() sort / reverse liste in place

Operations on Dictionaries

d[key]=value d.clear()
d[key] → value del d[key]
d.update(d2) update/add
d.keys() associations
d.values() iterable views on
d.items() keys/values/associations
d.pop(key[,default]) → value
d.popitem() → (key,value)
d.get(key[,default]) → value
d.setdefault(key[,default]) → value

Operations on Sets

Operators:
| → union (vertical bar char)
& → intersection
- ^ → difference/symmetric diff.
< <= > >= → inclusion relations
Operators also exist as methods.

s.update(s2) s.copy()
s.add(key) s.remove(key)
s.discard(key) s.clear()
s.pop()

Files

storing data on disk, and reading it back

f = open("file.txt", "w", encoding="utf8")

file variable name of file opening mode encoding of chars for text files:
for operations on disk (+path...) □ 'r' read □ 'w' write □ 'a' append utf8 ascii
cf. modules os, os.path and pathlib ... '+' 'x' 'b' 't' latin1 ...

writing f.write("coucou") reading f.read(n) → next chars
f.writelines(list of lines) if n not specified, read up to end!
f.readlines([n]) → list of next lines f.readline() → next line

text mode t by default (read/write str), possible binary mode b (read/write bytes). Convert from/to required type!

f.close() dont forget to close the file after use!

f.flush() write cache reading f.truncate([size]) resize
reading/writing progress sequentially in the file, modifiable with: f.tell() → position f.seek(position[,origin])

Very common: opening with a guarded block (automatic closing) and reading loop on lines with open(...) as f: for line in f: # processing of line

Iterative Loop Statement

statements block executed for each item of a container or iterator

for var in sequence:

→ statements block

next → finish

Go over sequence's values

s = "Some text" initializations before the loop
cnt = 0 loop variable, assignment managed by for statement
for var in s:
 if c == "e":
 cnt = cnt + 1
 print("found", cnt, "e")

loop on dict/set ⇔ loop on keys sequences use slices to loop on a subset of a sequence

Algo: count number of e in the string.

Display

print("v=", 3, "cm : ", x, ", ", y+4)

items to display : literal values, variables, expressions

Input

print options:

- sep=" " items separator, default space
- end="\n" end of print, default new line
- file=sys.stdout print to file, default standard output

Generic Operations on Containers

len(c) → items count
min(c) max(c) sum(c) Note: For dictionaries and sets, these operations use keys.
sorted(c) → list sorted copy
val in c → boolean, membership operator in (absence not in)
enumerate(c) → iterator on (index, value)
zip(c1, c2...) → iterator on tuples containing c1 items at same index
all(c) → True if all c items evaluated to true, else False
any(c) → True if at least one item of c evaluated true, else False
Specific to ordered sequences containers (lists, tuples, strings, bytes...)
reversed(c) → inverted iterator c*5 → duplicate c+c2 → concatenate
c.index(val) → position c.count(val) → events count
import copy
copy.copy(c) → shallow copy of container
copy.deepcopy(c) → deep copy of container

Operations on Lists

modify original list

lst.append(val) add item at end
lst.extend(seq) add sequence of items at end
lst.insert(idx, val) insert item at index
lst.remove(val) remove first item with value val
lst.pop(idx) → value remove & return item at index idx (default last)
lst.sort() lst.reverse() sort / reverse liste in place

Operations on Dictionaries

d[key]=value d.clear()
d[key] → value del d[key]
d.update(d2) update/add
d.keys() associations
d.values() iterable views on
d.items() keys/values/associations
d.pop(key[,default]) → value
d.popitem() → (key,value)
d.get(key[,default]) → value
d.setdefault(key[,default]) → value

Operations on Sets

Operators:
| → union (vertical bar char)
& → intersection
- ^ → difference/symmetric diff.
< <= > >= → inclusion relations
Operators also exist as methods.

s.update(s2) s.copy()
s.add(key) s.remove(key)
s.discard(key) s.clear()
s.pop()

Function Definition

function name (identifier) named parameters
def fct(x, y, z):
 """documentation"""
 # statements block, res computation, etc.
 return res ← result value of the call, if no computed
parameters and all result to return: return None
variables of this block exist only in the block and during the function call (think of a "black box")
Advanced: def fct(x, y, z, *args, a=3, b=5, **kwargs):
 *args variable positional arguments (→ tuple), default values,
 **kwargs variable named arguments (→ dict)

Function Call

r = fct(3, i+2, 2*i) storage/use of one argument per returned value
this is the use of function Advanced: *sequence **dict which does the call

Operations on Strings

s.startswith(prefix[,start[,end]]) s.endswith(suffix[,start[,end]]) s.strip(chars)
s.count(sub[,start[,end]]) s.partition(sep) → (before, sep, after)
s.index(sub[,start[,end]]) s.find(sub[,start[,end]])
s.is...() tests on chars categories (ex. s.isalpha())
s.upper() s.lower() s.title() s.swapcase()
s.casefold() s.capitalize() s.center(width,fill)
s.ljust(width,fill) s.rjust(width,fill) s.zfill(width)
s.encode(encoding) s.split(sep) s.join(seq)

Formatting

formatting directives values to format
"modele{} {} {}".format(x, y, r) → str
"selection:formatting!conversion"

Selection:
2 nom 0.nom 4[key] 0[2]

Formatting:
fill char alignment sign mini width.precision~maxwidth type
<> ^ + - space 0 at start for filling with 0
integer: b binary, c char, d decimal (default), o octal, x or X hexa...
float: e or E exponential, f or F fixed point, g or G appropriate (default), string: s ... % percent
Conversion: s (readable text) or x (literal representation)

Data Science Cheat Sheet

Python - Intermediate

KEY BASICS, PRINTING AND GETTING HELP

This cheat sheet assumes you are familiar with the content of our Python Basics Cheat Sheet

s - A Python string variable
i - A Python integer variable
f - A Python float variable

l - A Python list variable
d - A Python dictionary variable

LISTS

1.pop(3) - Returns the fourth item from **l** and deletes it from the list
1.remove(x) - Removes the first item in **l** that is equal to **x**
1.reverse() - Reverses the order of the items in **l**
l[1::2] - Returns every second item from **l**, commencing from the 1st item
l[-5:] - Returns the last 5 items from **l** specific axis

STRINGS

s.lower() - Returns a lowercase version of **s**
s.title() - Returns **s** with the first letter of every word capitalized
"23".zfill(4) - Returns "0023" by left-filling the string with 0's to make it's length 4.
s.splitlines() - Returns a list by splitting the string on any newline characters.
Python strings share some common methods with lists
s[:5] - Returns the first 5 characters of **s**
"fri" + "end" - Returns "friend"
"end" in s - Returns True if the substring "end" is found in **s**

RANGE

Range objects are useful for creating sequences of integers for looping.
range(5) - Returns a sequence from 0 to 4
range(2000,2018) - Returns a sequence from 2000 to 2017
range(0,11,2) - Returns a sequence from 0 to 10, with each item incrementing by 2
range(0,-10,-1) - Returns a sequence from 0 to -9
list(range(5)) - Returns a list from 0 to 4

DICTIONARIES

max(d, key=d.get) - Return the key that corresponds to the largest value in **d**
min(d, key=d.get) - Return the key that corresponds to the smallest value in **d**

SETS

my_set = set(l) - Return a **set** object containing the unique values from **l**

len(my_set) - Returns the number of objects in **my_set** (or, the number of unique values from **l**)
a in my_set - Returns True if the value **a** exists in **my_set**

REGULAR EXPRESSIONS

import re - Import the Regular Expressions module
re.search("abc",s) - Returns a **match** object if the regex "abc" is found in **s**, otherwise **None**
re.sub("abc","xyz",s) - Returns a string where all instances matching regex "abc" are replaced by "xyz"

LIST COMPREHENSION

A one-line expression of a for loop
[i ** 2 for i in range(10)] - Returns a list of the squares of values from 0 to 9
[s.lower() for s in l_strings] - Returns the list **l_strings**, with each item having had the **.lower()** method applied
[i for i in l_floats if i < 0.5] - Returns the items from **l_floats** that are less than 0.5

FUNCTIONS FOR LOOPING

```
for i, value in enumerate(l):
    print("The value of item {} is {}".
          format(i,value))
    - Iterate over the list l, printing the index location of each item and its value
```

```
for one, two in zip(l_one,l_two):
    print("one: {}, two: {}".format(one,two))
    - Iterate over two lists, l_one and l_two and print each value
```

```
while x < 10:
    x += 1
    - Run the code in the body of the loop until the value of x is no longer less than 10
```

DATETIME

import datetime as dt - Import the **datetime** module
now = dt.datetime.now() - Assign **datetime** object representing the current time to **now**
wks4 = dt.timedelta(weeks=4)
 - Assign a **timedelta** object representing a timespan of 4 weeks to **wks4**

now - wks4 - Return a **datetime** object representing the time 4 weeks prior to **now**

newyear_2020 = dt.datetime(year=2020, month=12, day=31) - Assign a **datetime** object representing December 25, 2020 to **newyear_2020**
newyear_2020.strftime("%A, %b %d, %Y")
 - Returns "Thursday, Dec 31, 2020"
dt.datetime.strptime('Dec 31, 2020', "%b %d, %Y") - Return a **datetime** object representing December 31, 2020

RANDOM

import random - Import the **random** module
random.random() - Returns a random float between 0.0 and 1.0
random.randint(0,10) - Returns a random integer between 0 and 10
random.choice(l) - Returns a random item from the list **l**

COUNTER

from collections import Counter - Import the **Counter** class
c = Counter(1) - Assign a **Counter** (dict-like) object with the counts of each unique item from 1, to **c**
c.most_common(3) - Return the 3 most common items from 1

TRY/EXCEPT

Catch and deal with Errors
l_ints = [1, 2, 3, "", 5] - Assign a list of integers with one missing value to **l_ints**
l_floats = []
for i in l_ints:
 try:
 l_floats.append(float(i))
 except:
 l_floats.append(i)
 - Convert each value of **l_ints** to a float, catching and handling **ValueError: could not convert string to float**: where values are missing.

Data Science Cheat Sheet

Python Regular Expressions

SPECIAL CHARACTERS

- \|** | Matches the expression to its right at the start of a string. It matches every such instance before each **\n** in the string.
- \\$|** | Matches the expression to its left at the end of a string. It matches every such instance before each **\n** in the string.
- .|** | Matches any character except line terminators like **\n**.
- \|** Escapes special characters or denotes character classes.
- A|B|** | Matches expression **A** or **B**. If **A** is matched first, **B** is left untried.
- +|** | Greedily matches the expression to its left 1 or more times.
- *|** | Greedily matches the expression to its left 0 or more times.
- ?|** | Greedily matches the expression to its left 0 or 1 times. But if **?** is added to qualifiers (+, *, and ? itself) it will perform matches in a non-greedy manner.
- {m}|** | Matches the expression to its left **m** times, and not less.
- {m,n}|** | Matches the expression to its left **m** to **n** times, and not less.
- {m,n}?|** | Matches the expression to its left **m** times, and ignores **n**. See **?** above.

CHARACTER CLASSES

[A.K.A. SPECIAL SEQUENCES]

- \w|** Matches alphanumeric characters, which means **a-z**, **A-Z**, and **0-9**. It also matches the underscore, **_**.
- \d|** Matches digits, which means **0-9**.
- \D|** Matches any non-digits.
- \s|** Matches whitespace characters, which include the **\t**, **\n**, **\r**, and space characters.
- \S|** Matches non-whitespace characters.
- \b|** Matches the boundary (or empty string) at the start and end of a word, that is, between **\w** and **\W**.
- \B|** Matches where **\b** does not, that is, the boundary of **\w** characters.

\A| Matches the expression to its right at the absolute start of a string whether in single or multi-line mode.

\Z| Matches the expression to its left at the absolute end of a string whether in single or multi-line mode.

SETS

- []|** Contains a set of characters to match.
- [amk]|** Matches either **a**, **m**, or **k**. It does not match **amk**.
- [a-z]|** Matches any alphabet from **a** to **z**.
- [a\z]|** Matches **a**, **-**, or **z**. It matches **-** because \ escapes it.
- [a-]|** Matches **a** or **-**, because **-** is not being used to indicate a series of characters.
- [-a]|** As above, matches **a** or **-**.
- [a-z0-9]|** Matches characters from **a** to **z** and also from **0** to **9**.
- [(+*)]|** Special characters become literal inside a set, so this matches **(, +, *, and)**.
- [^ab5]|** Adding **^** excludes any character in the set. Here, it matches characters that are not **a**, **b**, or **5**.

GROUPS

- ()|** Matches the expression inside the parentheses and groups it.
- (?)|** Inside parentheses like this, **?** acts as an extension notation. Its meaning depends on the character immediately to its right.
- (?PAB)|** Matches the expression **AB**, and it can be accessed with the group name.
- (?aiLmsux)|** Here, **a**, **i**, **L**, **m**, **s**, **u**, and **x** are flags:
 - a** — Matches ASCII only
 - i** — Ignore case
 - L** — Locale dependent
 - m** — Multi-line
 - s** — Matches all
 - u** — Matches unicode
 - x** — Verbose

(?:A)| Matches the expression as represented by **A**, but unlike **(?PAB)**, it cannot be retrieved afterwards.

(?#...)| A comment. Contents are for us to read, not for matching.

A(?=B)| Lookahead assertion. This matches the expression **A** only if it is followed by **B**.

A(?!B)| Negative lookahead assertion. This matches the expression **A** only if it is not followed by **B**.

(?<=B)A| Positive lookbehind assertion.

This matches the expression **A** only if **B** is immediately to its left. This can only matched fixed length expressions.

(?<!B)A| Negative lookbehind assertion.

This matches the expression **A** only if **B** is not immediately to its left. This can only matched fixed length expressions.

(?P=name)| Matches the expression matched by an earlier group named "name".

(...)\\1| The number **1** corresponds to the first group to be matched. If we want to match more instances of the same expression, simply use its number instead of writing out the whole expression again. We can use from **1** up to **99** such groups and their corresponding numbers.

POPULAR PYTHON RE MODULE FUNCTIONS

- re.findall(A, B)|** Matches all instances of an expression **A** in a string **B** and returns them in a list.
- re.search(A, B)|** Matches the first instance of an expression **A** in a string **B**, and returns it as a re match object.
- re.split(A, B)|** Split a string **B** into a list using the delimiter **A**.
- re.sub(A, B, C)|** Replace **A** with **B** in the string **C**.

Data Wrangling

with pandas Cheat Sheet

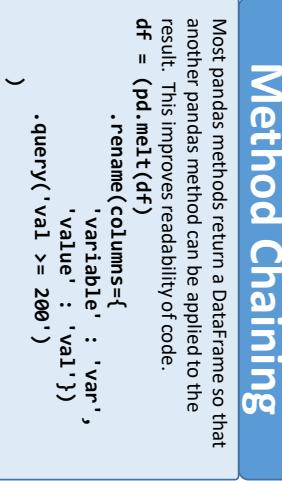
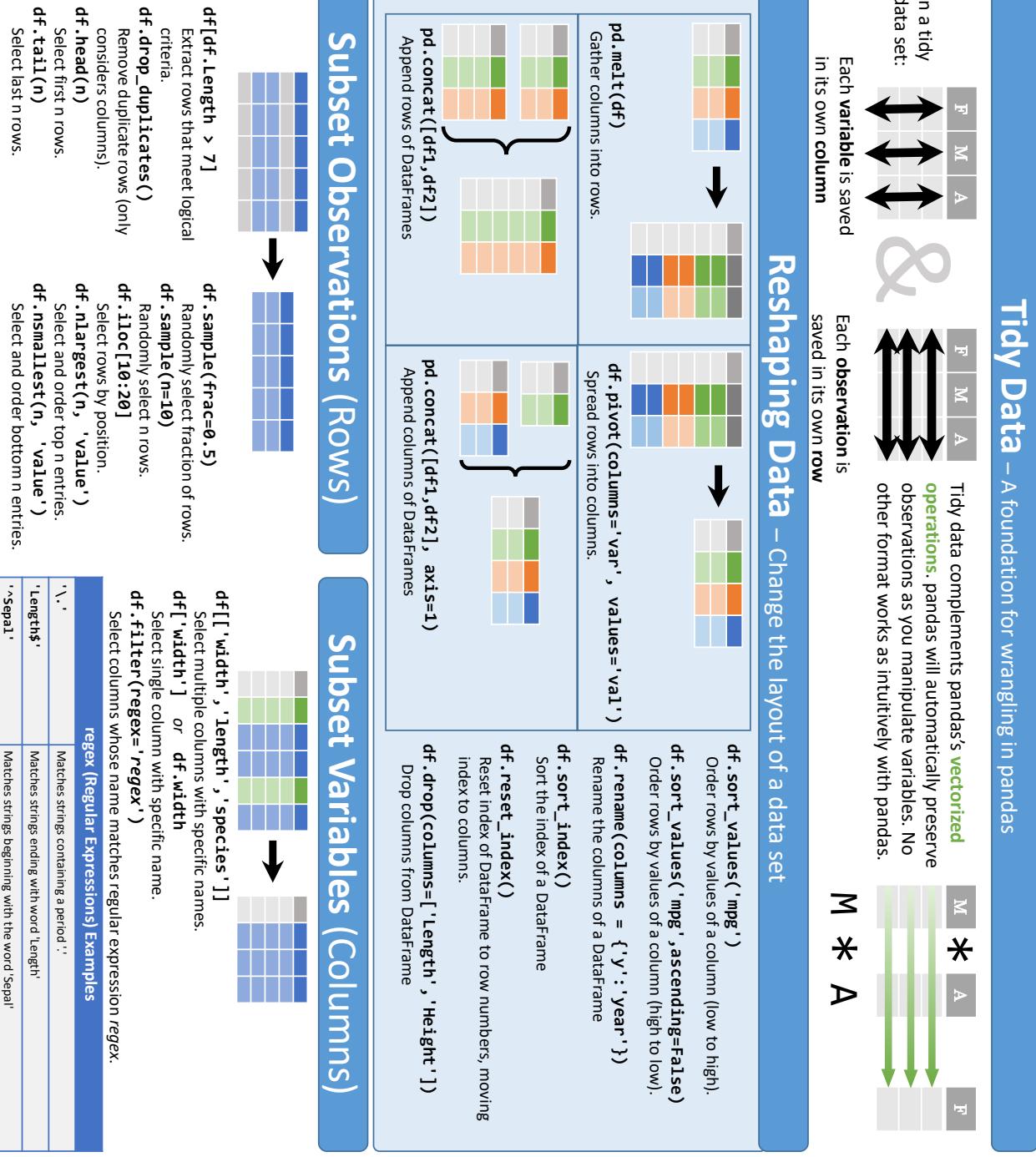
<http://pandas.pydata.org>

Syntax – Creating DataFrames

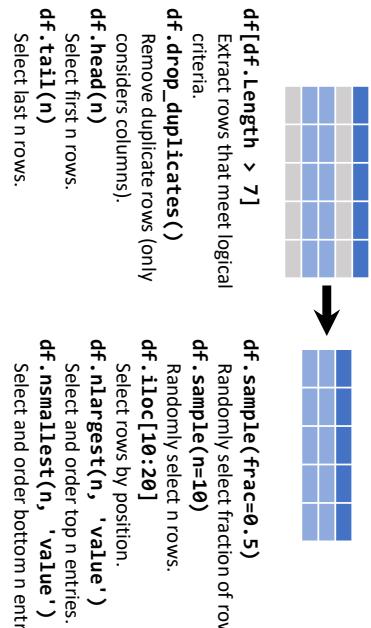
| | a | b | c |
|---|---|---|----|
| 1 | 4 | 7 | 10 |
| 2 | 5 | 8 | 11 |
| 3 | 6 | 9 | 12 |

```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.
```

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.
```



df = (pd.melt(df)
 .rename(columns={'variable': 'var',
 'value': 'val'})
 .query('val > 200')
)



Summarize Data

`df['w'].value_counts()`

Count number of rows with each unique value of variable

`len(df)`

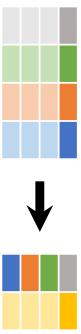
of rows in DataFrame.

`df['w'].nunique()`

of distinct values in a column.

`df.describe()`

Basic descriptive statistics for each column (or GroupBy)



pandas provides a large set of **summary functions** that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

`sum()`
Sum values of each object.
`count()`
Count non-NA/null values of each object.
`median()`
Median value of each object.
`Quantile([0.25, 0.75])`
Quantiles of each object.
`apply(function)`
Apply function to each object.

`min()`
Minimum value in each object.
`max()`
Maximum value in each object.
`mean()`
Mean value of each object.
`var()`
Variance of each object.
`std()`
Standard deviation of each object.

`min()`
Minimum value in each object.
`max()`
Maximum value in each object.
`pd.qcut(df.col, n, labels=False)`
Bin column into n buckets.

`df.assign(Area=lambda df: df.Length*df.Height)`
Compute and append one or more new columns.
`df['Volume'] = df.Length*df.Height*df.Depth`
Add single column.

`pd.qcut(df.col, n, labels=False)`
Bin column into n buckets.

`df.groupby(by='col')`
Return a GroupBy object, grouped by values in column named "col".
`df.groupby(level='ind')`
Return a GroupBy object, grouped by values in index level named "ind".

`apply(function)`
Apply function to each object.

`groupby('col')`
GroupBy object grouped by column.

`groupby('ind')`
GroupBy object grouped by index level.

`groupby(['col', 'ind'])`
GroupBy object grouped by both column and index level.

`groupby(['col', 'ind'], as_index=False)`
GroupBy object grouped by both column and index level, but do not make them columns.

`groupby(['col', 'ind'], group_keys=False)`
GroupBy object grouped by both column and index level, but do not make them keys.

`groupby(['col', 'ind'], dropna=True)`
GroupBy object grouped by both column and index level, but dropna=True.

`groupby(['col', 'ind'], dropna=False)`
GroupBy object grouped by both column and index level, but dropna=False.

`groupby(['col', 'ind'], group_keys=True)`
GroupBy object grouped by both column and index level, and make them keys.

`groupby(['col', 'ind'], group_keys=False)`
GroupBy object grouped by both column and index level, and do not make them keys.

`groupby(['col', 'ind'], as_index=True)`
GroupBy object grouped by both column and index level, but make them index.

`groupby(['col', 'ind'], as_index=False)`
GroupBy object grouped by both column and index level, but do not make them index.

`groupby(['col', 'ind'], as_index=None)`
GroupBy object grouped by both column and index level, but do not make them index.

`groupby(['col', 'ind'], as_index=True, dropna=True)`
GroupBy object grouped by both column and index level, but dropna=True.

`groupby(['col', 'ind'], as_index=True, dropna=False)`
GroupBy object grouped by both column and index level, but dropna=False.

`groupby(['col', 'ind'], as_index=False, dropna=True)`
GroupBy object grouped by both column and index level, but dropna=True.

`groupby(['col', 'ind'], as_index=False, dropna=False)`
GroupBy object grouped by both column and index level, but dropna=False.

`groupby(['col', 'ind'], as_index=None, dropna=True)`
GroupBy object grouped by both column and index level, but dropna=True.

`groupby(['col', 'ind'], as_index=None, dropna=False)`
GroupBy object grouped by both column and index level, but dropna=False.

`groupby(['col', 'ind'], as_index=None, dropna=None)`
GroupBy object grouped by both column and index level, but dropna=None.

Group Data

`df.groupby('col')`

Return a GroupBy object, grouped by values in column named "col".

`df.groupby(level='ind')`

Return a GroupBy object, grouped by values in index level named "ind".

`groupby('col')`

GroupBy object grouped by column.

`groupby('ind')`

GroupBy object grouped by index level.

`groupby(['col', 'ind'])`

GroupBy object grouped by both column and index level.

Windows

`df.plot.hist()`

Histogram for each column

`df.plot.scatter(x='w', y='h')`

Scatter chart using pairs of points

`df.expanding()`

Return an Expanding object allowing summary functions to be applied cumulatively.

`df.rolling(n)`

Return a Rolling object allowing summary functions to be applied to windows of length n.

Plotting

`df.plot.hist()`

Histogram for each column

`df.plot.scatter(x='w', y='h')`

Scatter chart using pairs of points

`df.exploding()`

Explode categorical variables in DataFrame.

`df.groupby(['col', 'ind']).size()`

Size of each group.

Combining Data Sets

`pd.merge(ydf, zdf, how='outer')`

Rows that appear in either or both ydf and zdf (Union).

`pd.merge(ydf, zdf, indicator=True)`

Indicator variable indicating which group each row belongs to.

`df.query('merge == "left_only"')`

Rows that appear in ydf but not zdf (Setdiff).

Handling Missing Data

`df.dropna()`

Drop rows with any column having NA/null data.

`df.fillna(value)`

Replace all NA/null data with value.

`df['w'].nunique()`

of distinct values in a column.

Make New Columns

`df.assign(Area=lambda df: df.Length*df.Height)`

Compute and append one or more new columns.

`df['Volume'] = df.Length*df.Height*df.Depth`

Add single column.

`pd.qcut(df.col, n, labels=False)`

Bin column into n buckets.

Standard Joins

`pd.merge(adf, bdf, how='left', on='x1')`

Join matching rows from bdf to adf.

`A 1 T`

`B 2 F`

`C 3 NaN`

`D NaN T`

Filtering Joins

`adf[~adf.x1.isin(bdf.x1)]`

All rows in adf that do not have a match in bdf.

`x1 x2`

`A 1`

`B 2`

Set-like Operations

`pd.merge(ydf, zdf, how='inner')`

Join data. Retain only rows in both sets.

`x1 x2 x3`

`A 1 T`

`B 2 F`

`C 3 NaN`

Combining Data Sets

`pd.merge(ydf, zdf, how='outer', indicator=True)`

Rows that appear in either or both ydf and zdf (Union).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='inner')`

Rows that appear in both ydf and zdf (Intersection).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='left')`

Rows that appear in ydf but not zdf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='right')`

Rows that appear in zdf but not ydf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='outer')`

Rows that appear in either ydf or zdf (Union).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='left_only')`

Rows that appear in ydf but not zdf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='right_only')`

Rows that appear in zdf but not ydf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='inner')`

Rows that appear in both ydf and zdf (Intersection).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='outer')`

Rows that appear in either ydf or zdf (Union).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='left')`

Rows that appear in ydf but not zdf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='right')`

Rows that appear in zdf but not ydf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='inner')`

Rows that appear in both ydf and zdf (Intersection).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='outer')`

Rows that appear in either ydf or zdf (Union).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='left_only')`

Rows that appear in ydf but not zdf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='right_only')`

Rows that appear in zdf but not ydf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='inner')`

Rows that appear in both ydf and zdf (Intersection).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='outer')`

Rows that appear in either ydf or zdf (Union).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='left')`

Rows that appear in ydf but not zdf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='right')`

Rows that appear in zdf but not ydf (Setdiff).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='inner')`

Rows that appear in both ydf and zdf (Intersection).

`x1 x2`

`A 1`

`B 2`

`C 3`

Plotting

`pd.merge(ydf, zdf, how='outer')`

Rows that appear in either ydf or zdf (Union).

`x1 x2`

`A 1`

`B 2`

<code

Python For Data Science Cheat Sheet

NumPy Basics

Learn Python for Data Science interactively at www.DataCamp.com

NumPy

The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:



NumPy Arrays

```
>>> import numpy as np
```

1D array

```
1 [ ] 2 [ ] 3 [ ]
```

axis 0 →

1.5 [] 2 [] 3 []

axis 1 →

4 [] 5 [] 6 []

axis 2 →

Python For Data Science Cheat Sheet

Excel Spreadsheets

Pickled Files

Learn Python for data science interactively at www.DataCamp.com



HDF5 Files

```
>>> import h5py  
>>> filename = 'H5_H1_LOSC_4_v1-815411200-4096.hdf5'  
>>> data = h5py.File(filename, 'r')
```

Matlab Files

```
>>> import scipy.io  
>>> filename = 'workspace.mat'  
>>> mat = scipy.io.loadmat(filename)
```

Exploring Dictionaries

| Accessing Data Items with Functions |
|---|
| >>> print(mat.keys()) >>> for key in mat.keys(): print(key) >>> meta = mat['meta'] >>> print(meta) >>> print(meta['Description']) >>> print(meta['URL']) >>> print(meta['Detector']) >>> print(meta['Duration']) >>> print(meta['GpsStart']) >>> print(meta['Observatory']) >>> print(meta['Type']) >>> print(meta['Urcaart']) >>> print(data['meta']['Description'].value) Retrieve the value for a key |

Navigating Your FileSystem

| OS Library |
|---|
| >>> import os >>> os.getcwd() List directory contents of files and directories >>> os.chdir('..') Change current working directory >>> os.getcwd() Return the current working directory path |

Magic Commands

| NumPy Arrays |
|--|
| >>> data = np.array([1, 2, 3]) Data type of array elements >>> data.dtype Array dimensions >>> len(data_array) Length of array |

Exploring Your Data

| pandas DataFrames |
|--|
| >>> df = pd.DataFrame({'A': [1, 2, 3], 'B': [4, 5, 6]}) >>> df.head(2) Return first DataFrame rows >>> df.tail(2) Return last DataFrame rows >>> df.columns Describe DataFrame columns >>> df.index Describe index >>> df.info() Info on DataFrame >>> data_array = df.values Convert DataFrame to a NumPy array |

Files with mixed data types

| CSV |
|---|
| >>> filename = 'titanic.csv' >>> data = np.loadtxt(filename, delimiter=',', skiprows=2, usecols=[0, 2], dtype=None) The default dtype of the np.loadtxt() function is None. |

Importing Flat Files with pandas

| CSV |
|---|
| >>> filename = 'winequality-red.csv' >>> data = pd.read_csv(filename, header=None, sep=';', comment='#', na_values=[''])) Number of rows offile to read Row number to use as col names Delimiter to use Character to split comments String to recognize as NA/NaN |

Importing Data in Python

Most of the time, you'll use either **NumPy** or **pandas** to import your data:

```
>>> import numpy as np  
>>> np.info(np.ndarray.dtype)  
>>> help(pd.read_csv)
```

Text Files

| Plain Text Files |
|--|
| >>> filename = 'huck_finn.txt' >>> file = open(filename, mode='r') Open the file for reading >>> text = file.read() Read a file's contents >>> print(file.closed) Check whether file is closed >>> file.close() >>> print(text) |

Using the context manager with
with open('huck_finn.txt', 'r') as file:
 print(file.readline())
print(file.readline())

Table Data: Flat Files

Importing Flat Files with numpy

Files with one data type

| CSV |
|---|
| >>> filename = 'huck_finn.txt' >>> data = np.loadtxt(filename, delimiter=',', skiprows=2, usecols=[0, 2], dtype=None) String used to separate values Skip the first 2 lines Read the 1st and 3rd column The type of the resulting array dtype=None) |

Importing Flat Files with pandas

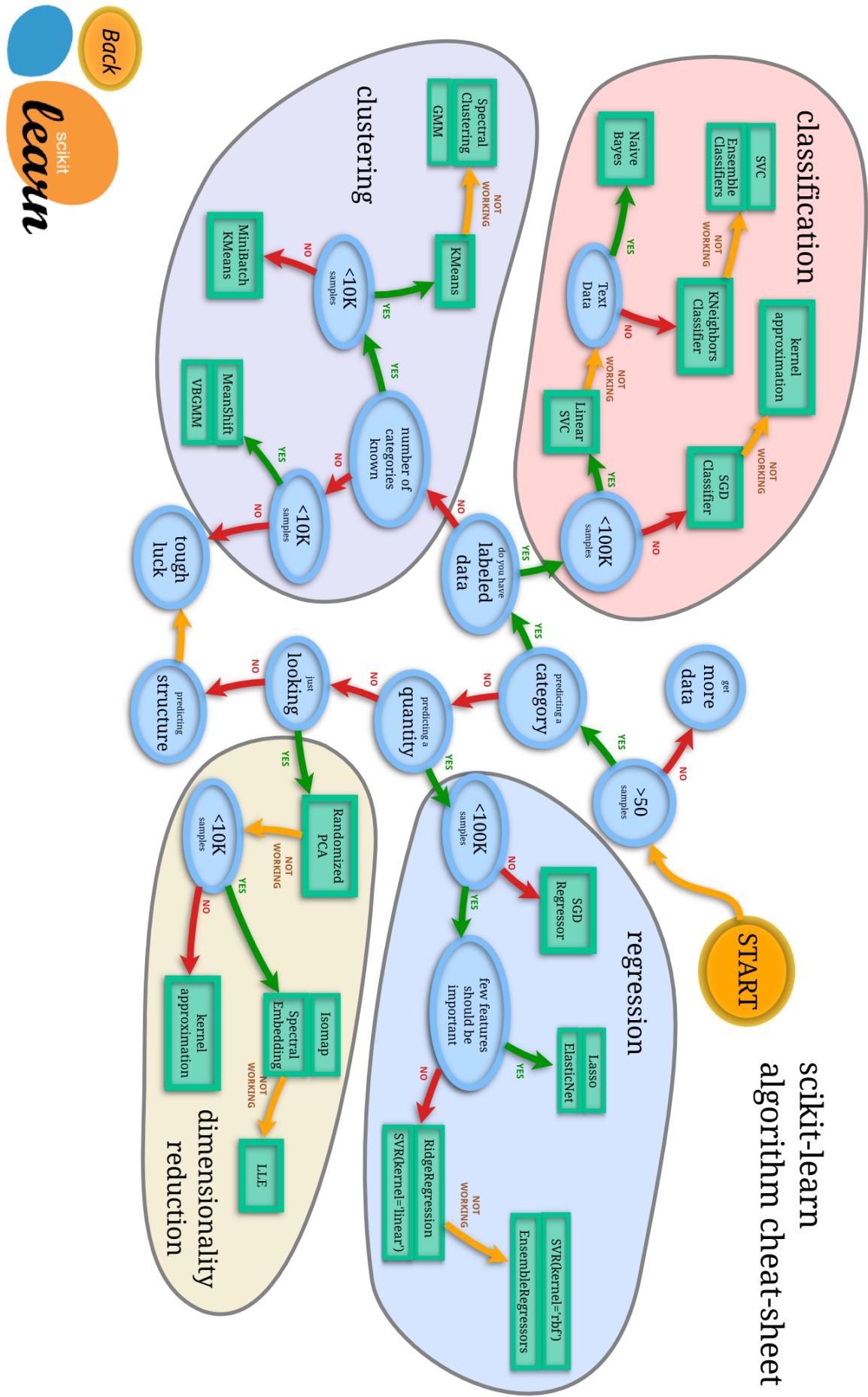
| CSV |
|---|
| >>> data = np.recfromcsv(filename) >>> data = np.genfromtxt(filename, delimiter=',', names=['text', 'id'], dtype=None) |

Importing Flat Files with pandas

| CSV |
|---|
| >>> data = np.recfromcsv(filename) >>> data = np.genfromtxt(filename, delimiter=',', names=['text', 'id'], dtype=None) |



scikit-learn algorithm cheat-sheet



Back
scikit
learn

Python For Data Science Cheat Sheet

Scikit-Learn

Scikit-learn is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.

A Basic Example

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.model_selection import train_test_split
>>> from sklearn.metrics import accuracy_score
>>> iris = datasets.load_iris()
>>> X, y = iris.data[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=33)
>>> X_train = preprocessing.StandardScaler().fit(X_train)
>>> X_test = scalar.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
```

KNN

```
>>> from sklearn.neighbors import KNeighborsClassifier
>>> gnb = GaussianNB()
>>> gnb = GaussianNB()
>>> gnb.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

Loading The Data

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as pandas DataFrame, are also acceptable.

```
>>> import numpy as np
>>> X = np.random.random((10, 5))
>>> y = np.array(['M', 'M', 'F', 'F', 'M', 'M', 'F', 'F', 'F'])
>>> X[X < 0.] = 0
```

Training And Test Data

```
>>> from sklearn.model_selection import train_test_split
>>> X_train, X_test, y_train, y_test = train_test_split(X, y,
... random_state=0)
```

Preprocessing The Data

Standardization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(X_train)
>>> standardized_X_train = scaler.transform(X_train)
```

Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X_train = scaler.transform(X_train)
>>> normalized_X_test = scaler.transform(X_test)
```

Binarization

```
>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0).fit(X)
>>> binary_X = binarizer.transform(X)
```

Create Your Model

Supervised Learning Estimators

Linear Regression

```
>>> from sklearn.linear_model import LinearRegression
>>> lr = LinearRegression(normalize=True)
```

Support Vector Machines (SVM)

```
>>> from sklearn.svm import SVC
```

Naive Bayes

```
>>> from sklearn.naive_bayes import GaussianNB
```

KNN

```
>>> from sklearn.neighbors import KNeighborsClassifier
```

Unsupervised Learning Estimators

Principal Component Analysis (PCA)

```
>>> from sklearn.decomposition import PCA
```

K-Means

```
>>> from sklearn.cluster import KMeans
```

Prediction

Supervised learning

```
>>> lr.fit(X, y)
```

Unsupervised Learning

```
>>> kmeans.fit(X_train)
```

Unsupervised Estimators

```
>>> pca_model = pca.fit_transform(X_train)
```

Fit the model to the data
Predict labels
Estimate probability of a label

Fit the model to the data
Predict labels
Estimate probability of a label

Predict labels in clustering algos

Model Fitting

Supervised Estimators

```
>>> lr.fit(X, y)
```

Unsupervised Estimators

```
>>> kmeans.fit(X_train)
```

Classification Metrics

Accuracy Score

```
>>> from sklearn.metrics import accuracy_score
```

Homogeneity

```
>>> from sklearn.metrics import homogeneity_score
```

V-measure

```
>>> from sklearn.metrics import v_measure_score
```

Cross-Validation

```
>>> from sklearn.cross_validation import cross_val_score
```

Tune Your Model

Grid Search

```
>>> from sklearn.grid_search import GridSearchCV
```

Adjusted Rand Index

```
>>> from sklearn.metrics import adjusted_rand_score
```

Mean Squared Error

```
>>> from sklearn.metrics import mean_squared_error
```

R2 Score

```
>>> from sklearn.metrics import r2_score
```

Clustering Metrics

Mean Absolute Error

```
>>> from sklearn.metrics import mean_absolute_error
```

Classification Report

```
>>> from sklearn.metrics import classification_report
```

Confusion Matrix

```
>>> from sklearn.metrics import confusion_matrix
```

Precision, recall, f1-score

```
>>> print(classification_report(y_test, y_pred))
```

Precision, recall, f1-score

```
>>> print(accuracy_score(y_true, y_pred))
```

Metric scoring functions

```
>>> accuracy_score(y_true, y_pred)
```

Precision, recall, f1-score

```
>>> print(precision_recall_fscore_support(y_true, y_pred))
```



VIP Cheatsheet: Machine Learning Tips

Afshine AMIDI and Shervine AMIDI

September 9, 2018

Metrics

Given a set of data points $\{x^{(1)}, \dots, x^{(m)}\}$, where each $x^{(i)}$ has n features, associated to a set of outcomes $\{y^{(1)}, \dots, y^{(m)}\}$, we want to assess a given classifier that learns how to predict y from x .

In a context of a binary classification, here are the main metrics that are important to track to assess the performance of the model.

□ **Confusion matrix** – The confusion matrix is used to have a more complete picture when assessing the performance of a model. It is defined as follows:

| | | Predicted class | |
|--------------|---|-----------------------|--|
| | | + | - |
| Actual class | + | TP True Positives | FN False Negatives Type II error |
| | - | FP False Positives | TN True Negatives |

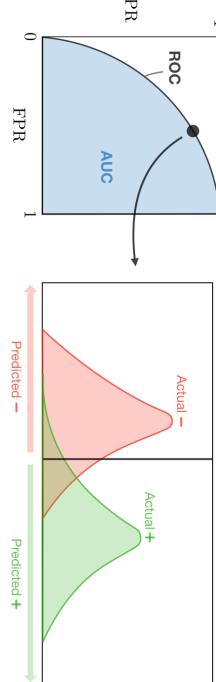
□ **Main metrics** – The following metrics are commonly used to assess the performance of classification models:

| Metric | Formula | Interpretation |
|-------------|---|---|
| Accuracy | $\frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$ | Overall performance of model |
| Precision | $\frac{\text{TP}}{\text{TP} + \text{FP}}$ | How accurate the positive predictions are |
| Recall | $\frac{\text{TP}}{\text{TP} + \text{FN}}$ | Coverage of actual positive sample |
| Sensitivity | $\frac{\text{TP}}{\text{TP} + \text{FN}}$ | Coverage of actual negative sample |
| Specificity | $\frac{\text{TN}}{\text{TN} + \text{FP}}$ | Coverage of actual negative sample |
| F1 score | $\frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$ | Hybrid metric useful for unbalanced classes |

□ **ROC** – The receiver operating curve, also noted ROC, is the plot of TPR versus FPR by varying the threshold. These metrics are summed up in the table below:

| Metric | Formula | Equivalent |
|---------------------|---|---------------------|
| True Positive Rate | $\frac{\text{TP}}{\text{TP} + \text{FN}}$ | Recall, sensitivity |
| False Positive Rate | $\frac{\text{FP}}{\text{TN} + \text{FP}}$ | 1-specificity |

□ **AUC** – The area under the receiving operating curve, also noted AUC or AUROC, is the area below the ROC as shown in the following figure:



Regression

□ **Basic metrics** – Given a regression model f , the following metrics are commonly used to assess the performance of the model:

| Total sum of squares | Explained sum of squares | Residual sum of squares |
|---|--|--|
| $\text{SS}_{\text{tot}} = \sum_{i=1}^m (y_i - \bar{y})^2$ | $\text{SS}_{\text{reg}} = \sum_{i=1}^m (f(x_i) - \bar{y})^2$ | $\text{SS}_{\text{res}} = \sum_{i=1}^m (y_i - f(x_i))^2$ |

□ **Coefficient of determination** – The coefficient of determination, often noted R^2 or r^2 , provides a measure of how well the observed outcomes are replicated by the model and is defined as follows:

$$R^2 = 1 - \frac{\text{SS}_{\text{res}}}{\text{SS}_{\text{tot}}}$$

□ **Main metrics** – The following metrics are commonly used to assess the performance of regression models, by taking into account the number of variables n that they take into consideration:

| Mallow's Cp | AIC | BIC | Adjusted R^2 |
|---|----------------------|---------------------------|----------------------------------|
| $\frac{\text{SS}_{\text{res}} + 2(n+1)\hat{\sigma}^2}{m}$ | $2[(n+2) - \log(L)]$ | $\log(m)(n+2) - 2\log(L)$ | $1 - \frac{(1-R^2)(m-1)}{m-n-1}$ |

where L is the likelihood and $\hat{\sigma}^2$ is an estimate of the variance associated with each response.

Model selection

Vocabulary – When selecting a model, we distinguish 3 different parts of the data that we have as follows:

| Training set | Validation set | Testing set |
|--|---|--|
| - Model is trained - Usually 80% of the dataset or development set | - Model is assessed - Usually 20% of the dataset - Also called hold-out | - Model gives predictions - Unseen data |

Once the model has been chosen, it is trained on the entire dataset and tested on the unseen test set. These are represented in the figure below:



Cross-validation – Cross-validation, also noted CV, is a method that is used to select a model that does not rely too much on the initial training set. The different types are summed up in the table below:

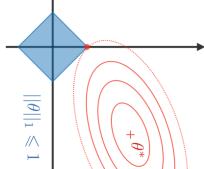
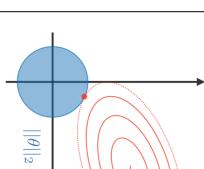
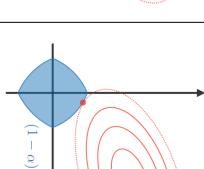
| k-fold | Leave-p-out |
|--|---|
| - Training on $k - 1$ folds and assessment on the remaining one - Generally $k = 5$ or 10 | - Training on $n - p$ observations and assessment on the p remaining ones - Case $p = 1$ is called leave-one-out |

The most commonly used method is called k -fold cross-validation and splits the training data into k folds to validate the model on one fold while training the model on the $k - 1$ other folds, all of this k times. The error is then averaged over the k folds and is named cross-validation error.

| Fold | Dataset | Validation error | Cross-validation error |
|------|----------------|------------------|---|
| 1 | (Green circle) | ϵ_1 | |
| 2 | (Red circle) | ϵ_2 | $\frac{\epsilon_1 + \dots + \epsilon_k}{k}$ |
| : | : | : | |
| k | (Green circle) | ϵ_k | |

Regularization – The regularization procedure aims at avoiding the model to overfit the data and thus deals with high variance issues. The following table sums up the different types of commonly used regularization techniques:

| LASSO | Ridge | Elastic Net |
|--|----------------------------|--|
| - Shrinks coefficients to 0 - Good for variable selection | Makes coefficients smaller | Tradeoff between variable selection and small coefficients |

| | | |
|---|---|---|
|  |  |  |
| $ \theta _1 \leqslant 1$ | $ \theta _2 \leqslant 1$ | $(1 - \alpha) \theta _1 + \alpha \theta _2^2 \leqslant 1$ |
| $\dots + \lambda \theta _1$ | $\dots + \lambda \theta _2$ | $\dots + \lambda \left[(1 - \alpha) \theta _1 + \alpha \theta _2^2 \right]$ |
| $\lambda \in \mathbb{R}$ | $\lambda \in \mathbb{R}$ | $\lambda \in \mathbb{R}, \alpha \in [0,1]$ |

Model selection – Train model on training set, then evaluate on the development set, then pick best performance model on the development set, and retrain all of that model on the whole training set.

Diagnostics

Bias – The bias of a model is the difference between the expected prediction and the correct model that we try to predict for given data points.

Variance – The variance of a model is the variability of the model prediction for given data points.

Bias/variance tradeoff – The simpler the model, the higher the bias, and the more complex the model, the higher the variance.

| Underfitting | Just right | Overfitting |
|--|---|--|
| - High training error - Training error close to test error - High bias | - Training error slightly lower than test error | - Low training error - Training error much lower than test error - High variance |

| | | | |
|----------------|---|--|---|
| Classification | | | |
| Deep learning | | | |
| Remedies | <ul style="list-style-type: none"> - Complexity model - Add more features - Train longer | | <ul style="list-style-type: none"> - Regularize - Get more data |

- **Error analysis** – Error analysis is analyzing the root cause of the difference in performance between the current and the perfect models.
- **Ablative analysis** – Ablative analysis is analyzing the root cause of the difference in performance between the current and the baseline models.

VIP Cheatsheet: Supervised Learning

Afshine AMIDI and Shervine AMIDI

September 9, 2018

Introduction to Supervised Learning

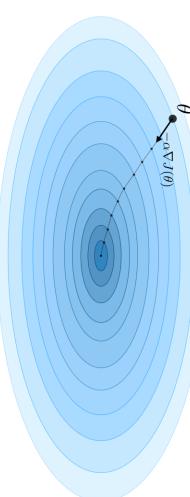
Given a set of data points $\{x^{(1)}, \dots, x^{(m)}\}$ associated to a set of outcomes $\{y^{(1)}, \dots, y^{(m)}\}$, we want to build a classifier that learns how to predict y from x .

□ **Type of prediction** – The different types of predictive models are summed up in the table below:

| Regression | Classifier |
|------------|-------------------------------------|
| Outcome | Continuous Class |
| Examples | Linear regression, SVM, Naive Bayes |

□ **Type of model** – The different models are summed up in the table below:

| Discriminative model | Generative model |
|----------------------|--|
| Goal | Directly estimate $P(y x)$ |
| What's learned | Decision boundary Probability distributions of the data |



Remark: Stochastic gradient descent (SGD) is updating the parameter based on each training example, and batch gradient descent is on a batch of training examples.

□ **Likelihood** – The likelihood of a model $L(\theta)$ given parameters θ is used to find the optimal parameters θ through maximizing the likelihood. In practice, we use the log-likelihood $\ell(\theta) = \log(L(\theta))$ which is easier to optimize. We have:

$$\theta^{\text{opt}} = \arg \max_{\theta} L(\theta)$$

□ **Newton's algorithm** – The Newton's algorithm is a numerical method that finds θ such that $\ell'(\theta) = 0$. Its update rule is as follows:

$$\theta \leftarrow \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

□ **Hypothesis** – The hypothesis is noted h_{θ} and is the model that we choose. For a given input data $x^{(i)}$, the model prediction output is $h_{\theta}(x^{(i)})$.

□ **Loss function** – A loss function is a function $L : (z, y) \in \mathbb{R} \times Y \mapsto L(z, y) \in \mathbb{R}$ that takes as inputs the predicted value z corresponding to the real data value y and outputs how different they are. The common loss functions are summed up in the table below:

| Least squared | Logistic | Hinge | Cross-entropy |
|------------------------|-----------------------|-------------------|---|
| $\frac{1}{2}(y - z)^2$ | $\log(1 + \exp(-yz))$ | $\max(0, 1 - yz)$ | $-\left[y \log(z) + (1 - y) \log(1 - z)\right]$ |
| | | | |
| Linear regression | Logistic regression | SVM | Neural Network |

□ **Cost function** – The cost function J is commonly used to assess the performance of a model, and is defined with the loss function L as follows:

$$J(\theta) = \sum_{i=1}^m L(h_{\theta}(x^{(i)}), y^{(i)})$$

□ **Gradient descent** – By noting $\alpha \in \mathbb{R}$ the learning rate, the update rule for gradient descent is expressed with the learning rate and the cost function J as follows:

$$\theta \leftarrow \theta - \alpha \nabla J(\theta)$$

Linear regression

We assume here that $y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$

\square **Normal equations** – By noting X the matrix design, the value of θ that minimizes the cost function is a closed-form solution such that:

$$\theta = (X^T X)^{-1} X^T y$$

\square **LMS algorithm** – By noting α the learning rate, the update rule of the Least Mean Squares (LMS) algorithm for a training set of m data points, which is also known as the Widrow-Hoff learning rule, is as follows:

$$\forall j, \quad \theta_j \leftarrow \theta_j + \alpha \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] x_j^{(i)}$$

Remark: the update rule is a particular case of the gradient ascent.

\square **LWR** – Locally Weighted Regression, also known as LWR, is a variant of linear regression that weights each training example in its cost function by $w^{(i)}(x)$, which is defined with parameter $\tau \in \mathbb{R}$ as:

$$w^{(i)}(x) = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$

Classification and logistic regression

\square **Sigmoid function** – The sigmoid function g , also known as the logistic function, is defined as follows:

$$\forall z \in \mathbb{R}, \quad g(z) = \frac{1}{1 + e^{-z}} \in [0, 1]$$

\square **Assumptions of GLMs** – Generalized Linear Models (GLM) aim at predicting a random variable y as a function to $x \in \mathbb{R}^{n+1}$ and rely on the following 3 assumptions:

- (1) $y|x; \theta \sim \text{ExpFamily}(\eta)$
- (2) $h_\theta(x) = E[y|x; \theta]$
- (3) $\eta = \theta^T x$

Remark: ordinary least squares and logistic regression are special cases of generalized linear models.

\square **Logistic regression** – We assume here that $y|x; \theta \sim \text{Bernoulli}(\phi)$. We have the following form:

$$\phi = p(y = 1|x; \theta) = \frac{1}{1 + \exp(-\theta^T x)} = g(\theta^T x)$$

Remark: there is no closed form solution for the case of logistic regressions.

\square **Softmax regression** – A softmax regression, also called a multiclass logistic regression, is used to generalize logistic regression when there are more than 2 outcome classes. By convention, we set $\theta_K = 0$, which makes the Bernoulli parameter ϕ_i of each class i equal to:

$$\phi_i = \frac{\exp(\theta_i^T x)}{\sum_{j=1}^K \exp(\theta_j^T x)}$$

where $(w, b) \in \mathbb{R}^n \times \mathbb{R}$ is the solution of the following optimization problem:

$$\min \frac{1}{2} \|w\|^2 \quad \text{such that} \quad y^{(i)}(w^T x^{(i)} - b) \geq 1$$

Generalized Linear Models

\square **Exponential family** – A class of distributions is said to be in the exponential family if it can be written in terms of a natural parameter, also called the canonical parameter or link function, η , a sufficient statistic $T(y)$ and a log-partition function $a(\eta)$ as follows:

$$p(y; \eta) = b(y) \exp(\eta T(y) - a(\eta))$$

Remark: we will often have $T(y) = y$. Also, $\exp(-a(\eta))$ can be seen as a normalization parameter that will make sure that the probabilities sum to one.

Here are the most common exponential distributions summed up in the following table:

| Distribution | η | $T(y)$ | $a(\eta)$ | $b(y)$ |
|--------------|--|--------|--|---|
| Bernoulli | $\log\left(\frac{\phi}{1-\phi}\right)$ | y | $\log(1 + \exp(\eta))$ | 1 |
| Gaussian | μ | y | $\frac{\eta^2}{2}$ | $\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$ |
| Poisson | $\log(\lambda)$ | y | e^η | $\frac{1}{y!}$ |
| Geometric | $\log(1 - \phi)$ | y | $\log\left(\frac{e^\eta}{1-e^\eta}\right)$ | 1 |

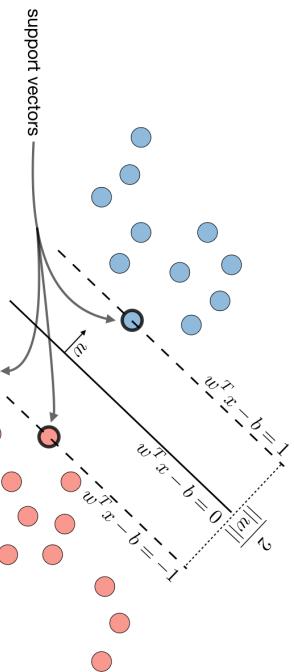
Gaussian Discriminant Analysis

□ **Setting** – The Gaussian Discriminant Analysis assumes that y and $x|y = 0$ and $x|y = 1$ are such that:

$$y \sim \text{Bernoulli}(\phi)$$

$$x|y = 0 \sim \mathcal{N}(\mu_0, \Sigma)$$

$$x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)$$



Remark: the line is defined as $w^T x - b = 0$.

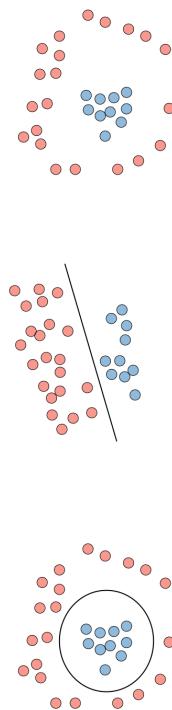
□ **Hinge loss** – The hinge loss is used in the setting of SVMs and is defined as follows:

$$L(z,y) = [1 - yz]_+ = \max(0, 1 - yz)$$

□ **Kernel** – Given a feature mapping ϕ , we define the kernel K to be defined as:

$$K(x,z) = \phi(x)^T \phi(z)$$

In practice, the kernel K defined by $K(x,z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$ is called the Gaussian kernel and is commonly used.



Non-linear separability → Use of a kernel mapping ϕ → Decision boundary in the original space

Remark: we say that we use the "kernel trick" to compute the cost function using the kernel because we actually don't need to know the explicit mapping ϕ , which is often very complicated. Instead, only the values $K(x,z)$ are needed.

□ **Lagrangian** – We define the Lagrangian $\mathcal{L}(w,b)$ as follows:

$$\mathcal{L}(w,b) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

Remark: the coefficients β_i are called the Lagrange multipliers.

Generative Learning

A generative model first tries to learn how the data is generated by estimating $P(x|y)$, which we can then use to estimate $P(y|x)$ by using Bayes' rule.

Naive Bayes

□ **Solutions** – Maximizing the log-likelihood gives the following solutions, with $k \in \{0,1\}$, $l \in [1, L]$

$$P(x|y) = P(x_1, x_2, \dots | y) = P(x_1|y)P(x_2|y)\dots = \prod_{i=1}^n P(x_i|y)$$

Remark: Naive Bayes is widely used for text classification and spam detection.

Tree-based and ensemble methods

These methods can be used for both regression and classification problems.

□ **CART** – Classification and Regression Trees (CART), commonly known as decision trees, can be represented as binary trees. They have the advantage to be very interpretable.

□ **Random forest** – It is a tree-based technique that uses a high number of decision trees built out of randomly selected sets of features. Contrary to the simple decision tree, it is highly uninterpretable but its generally good performance makes it a popular algorithm.

Remark: random forests are a type of ensemble methods.

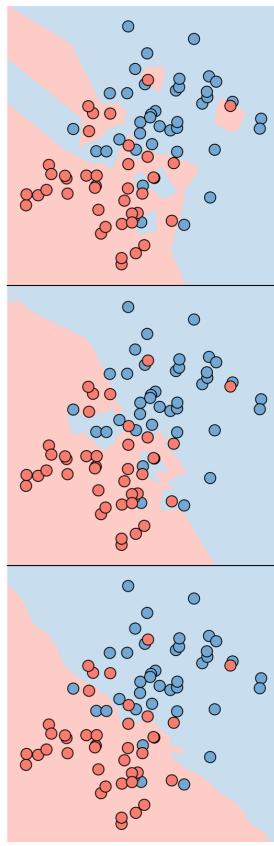
□ **Boosting** – The idea of boosting methods is to combine several weak learners to form a stronger one. The main ones are summed up in the table below:

| Adaptive boosting | Gradient boosting |
|--|---|
| <ul style="list-style-type: none"> - High weights are put on errors to improve at the next boosting step - Known as AdaBoost | <ul style="list-style-type: none"> - Weak learners trained on remaining errors |

Other non-parametric approaches

□ **k -nearest neighbors** – The k -nearest neighbors algorithm, commonly known as k -NN, is a non-parametric approach where the response of a data point is determined by the nature of its k neighbors from the training set. It can be used in both classification and regression settings.

Remark: The higher the parameter k , the higher the bias, and the lower the parameter k , the higher the variance.



$k = 1$

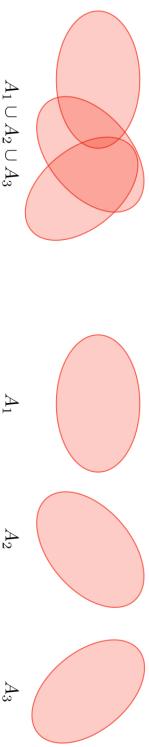
$k = 3$

$k = 11$

Learning Theory

□ **Union bound** – Let A_1, \dots, A_k be k events. We have:

$$P(A_1 \cup \dots \cup A_k) \leq P(A_1) + \dots + P(A_k)$$



$A_1 \cup A_2 \cup A_3$

□ **Hoeffding inequality** – Let Z_1, \dots, Z_m be m iid variables drawn from a Bernoulli distribution of parameter ϕ . Let $\hat{\phi}$ be their sample mean and $\gamma > 0$ fixed. We have:

$$P(|\phi - \hat{\phi}| > \gamma) \leq 2 \exp(-2\gamma^2 m)$$

Remark: this inequality is also known as the Chernoff bound.

□ **Training error** – For a given classifier h , we define the training error $\hat{e}(h)$, also known as the empirical risk or empirical error, to be as follows:

$$\hat{e}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x^{(i)}) \neq y^{(i)}\}}$$

□ **Probably Approximately Correct (PAC)** – PAC is a framework under which numerous results on learning theory were proved, and has the following set of assumptions:

- the training and testing sets follow the same distribution
- the training examples are drawn independently

□ **Shattering** – Given a set $S = \{x^{(1)}, \dots, x^{(d)}\}$, and a set of classifiers \mathcal{H} , we say that \mathcal{H} shatters S if for any set of labels $\{y^{(1)}, \dots, y^{(d)}\}$, we have:

$$\exists h \in \mathcal{H}, \quad \forall i \in [1, d], \quad h(x^{(i)}) = y^{(i)}$$

□ **VC dimension** – The Vapnik-Chervonenkis (VC) dimension of a given infinite hypothesis class \mathcal{H} , noted $\text{VC}(\mathcal{H})$ is the size of the largest set that is shattered by \mathcal{H} .

Remark: the VC dimension of $\mathcal{H} = \{\text{set of linear classifiers in 2 dimensions}\}$ is 3.



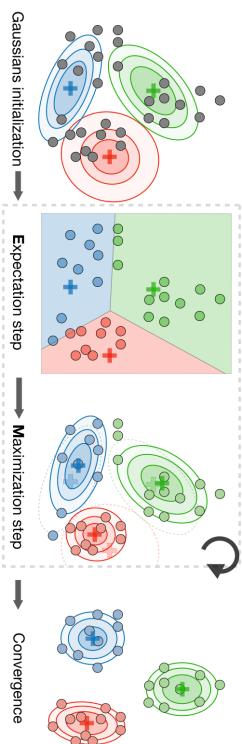
□ **Theorem (Vapnik)** – Let \mathcal{H} be given, with $\text{VC}(\mathcal{H}) = d$ and m the number of training examples. With probability at least $1 - \delta$, we have:

$$\epsilon(\hat{h}) \leq \left(\min_{h \in \mathcal{H}} \epsilon(h) \right) + O\left(\sqrt{\frac{d}{m} \log\left(\frac{m}{d}\right)} + \frac{1}{m} \log\left(\frac{1}{\delta}\right) \right)$$

VIP Cheatsheet: Unsupervised Learning

Afshine AMIDI and Shervine AMIDI

September 9, 2018



Introduction to Unsupervised Learning

Motivation – The goal of unsupervised learning is to find hidden patterns in unlabeled data $\{x^{(1)}, \dots, x^{(m)}\}$.

Jensen's inequality – Let f be a convex function and X a random variable. We have the following inequality:

$$E[f(X)] \geq f(E[X])$$

Expectation-Maximization

Latent variables – Latent variables are hidden/unobserved variables that make estimation problems difficult, and are often denoted z . Here are the most common settings where there are latent variables:

| Setting | Latent variable z | $x z$ | Comments |
|--------------------------|----------------------------|------------------------------------|---|
| Mixture of k Gaussians | $\text{Multinomial}(\phi)$ | $\mathcal{N}(\mu_j, \Sigma_j)$ | $\mu_j \in \mathbb{R}^n, \phi \in \mathbb{R}^k$ |
| Factor analysis | $\mathcal{N}(0, I)$ | $\mathcal{N}(\mu + \Lambda z\psi)$ | $\mu_j \in \mathbb{R}^n$ |

Algorithm – The Expectation-Maximization (EM) algorithm gives an efficient method at estimating the parameter θ through maximum likelihood estimation by repeatedly constructing a lower-bound on the likelihood (E-step) and optimizing that lower bound (M-step) as follows:

- **E-step:** Evaluate the posterior probability $Q_i(z^{(i)})$ that each data point $x^{(i)}$ came from a particular cluster $z^{(i)}$ as follows:

$$Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

- **M-step:** Use the posterior probabilities $Q_i(z^{(i)})$ as cluster specific weights on data points $x^{(i)}$ to separately re-estimate each cluster model as follows:

$$\theta_i = \operatorname{argmax}_{\theta} \sum_i \int_{z^{(i)}} Q_i(z^{(i)}) \log \left(\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) dz^{(i)}$$

Hierarchical clustering

Algorithm – It is a clustering algorithm with an agglomerative hierarchical approach that build nested clusters in a successive manner.

Types – There are different sorts of hierarchical clustering algorithms that aims at optimizing different objective functions, which is summed up in the table below:

| Ward linkage | Average linkage | Complete linkage |
|----------------------------------|---|--|
| Minimize within cluster distance | Minimize average distance between cluster pairs | Minimize maximum distance of between cluster pairs |

Clustering assessment metrics

In an unsupervised learning setting, it is often hard to assess the performance of a model since we don't have the ground truth labels as was the case in the supervised learning setting.

□ **Silhouette coefficient** – By noting a and b the mean distance between a sample and all other points in the same class, and between a sample and all other points in the next nearest cluster, the silhouette coefficient s for a single sample is defined as follows:

$$s = \frac{b - a}{\max(a, b)}$$

□ **Calinski-Harabaz index** – By noting k the number of clusters, B_k and W_k the between and within-clustering dispersion matrices respectively defined as

$$B_k = \sum_{j=1}^k n_{c^{(j)}} (\mu_{c^{(j)}} - \mu) (\mu_{c^{(j)}} - \mu)^T, \quad W_k = \sum_{i=1}^m (x^{(i)} - \mu_{c^{(i)}}) (x^{(i)} - \mu_{c^{(i)}})^T$$

the Calinski-Harabaz index $s(k)$ indicates how well a clustering model defines its clusters, such that the higher the score, the more dense and well separated the clusters are. It is defined as follows:

$$s(k) = \frac{\text{Tr}(B_k)}{\text{Tr}(W_k)} \times \frac{N - k}{k - 1}$$

Principal component analysis

It is a dimension reduction technique that finds the variance maximizing directions onto which to project the data.

□ **Eigenvalue, eigenvector** – Given a matrix $A \in \mathbb{R}^{n \times n}$, λ is said to be an eigenvalue of A if there exists a vector $z \in \mathbb{R}^n \setminus \{0\}$, called eigenvector, such that we have:

$$Az = \lambda z$$

□ **Spectral theorem** – Let $A \in \mathbb{R}^{n \times n}$. If A is symmetric, then A is diagonalizable by a real orthogonal matrix $U \in \mathbb{R}^{n \times n}$. By noting $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$, we have:

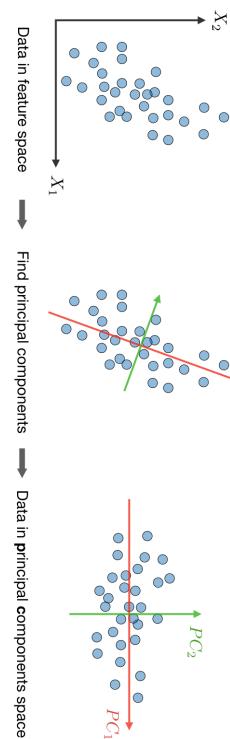
$$\exists \Lambda \text{ diagonal}, \quad A = U \Lambda U^T$$

Remark: the eigenvector associated with the largest eigenvalue is called principal eigenvector of matrix A .

□ **Algorithm** – The Principal Component Analysis (PCA) procedure is a dimension reduction technique that projects the data on k dimensions by maximizing the variance of the data as follows:

- Step 1: Normalize the data to have a mean of 0 and standard deviation of 1.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j} \quad \text{where} \quad \mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{and} \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$



Independent component analysis

It is a technique meant to find the underlying generating sources.

□ **Assumptions** – We assume that our data x has been generated by the n -dimensional source vector $s = (s_1, \dots, s_n)$, where s_i are independent random variables, via a mixing and non-singular matrix A as follows:

$$x = As$$

The goal is to find the unmixing matrix $W = A^{-1}$ by an update rule.

□ **Bell and Sejnowski ICA algorithm** – This algorithm finds the unmixing matrix W by following the steps below:

- Write the probability of $x = As = W^{-1}s$ as:

$$p(x) = \prod_{i=1}^n p_s(w_i^T x) \cdot |W|$$

- Write the log likelihood given our training data $\{x^{(i)}, i \in [1, m]\}$ and by noting g the sigmoid function as:

$$l(W) = \sum_{i=1}^m \left(\sum_{j=1}^n \log \left(g'(w_j^T x^{(i)}) \right) + \log |W| \right)$$

Therefore, the stochastic gradient ascent learning rule is such that for each training example $x^{(i)}$, we update W as follows:

$$W \leftarrow W + \alpha \begin{pmatrix} \frac{1 - 2g(w_1^T x^{(i)})}{1 - 2g(w_2^T x^{(i)})} \\ \vdots \\ \frac{1 - 2g(w_n^T x^{(i)})}{1 - 2g(w_1^T x^{(i)})} \end{pmatrix} x^{(i)T} + (W^T)^{-1}$$

Python For Data Science Cheat Sheet

Model Architecture

Sequential Model

Learn Python for data science interactively at www.DataCamp.com



Keras

Keras is a powerful and easy-to-use deep learning library for Theano and TensorFlow that provides a high-level neural networks API to develop and evaluate deep learning models.

A Basic Example

```
>>> import numpy as np
>>> from keras.models import Sequential
>>> from keras.layers import Dense
>>> data = np.random.randint(2, size=(1000, 1))
>>> model = Sequential()
>>> model.add(Dense(32,
>>>                 activation='relu',
>>>                 input_dim=100))
>>> model.add(Dense(1, activation='sigmoid'))
>>> model.compile(optimizer='rmsprop',
>>>                 loss='binary_crossentropy',
>>>                 metrics=['accuracy'])
>>> model.fit(data, labels, epochs=10, batch_size=32)
>>> predictions = model.predict(data)
```

Data

Your data needs to be stored as NumPy arrays or as a list of NumPy arrays. Ideally, you'll split the data in training and test sets, for which you can also resort to the `train_test_split` module of `sklearn.cross_validation`.

Keras Data Sets

```
>>> from keras.datasets import boston_housing,
>>>                      mnist,
>>>                      imdb,
>>>                      cifar10,
>>>                      (x_train, y_train), (x_test, y_test) = boston_housing.load_data()
>>> (x_train2, y_train2), (x_test2, y_test2) = mnist.load_data()
>>> (x_train3, y_train3), (x_test3, y_test3) = cifar10.load_data()
>>> (x_train4, y_train4), (x_test4, y_test4) = imdb.load_data(maxlen=20000)
>>> num_classes = 10
```

Other

```
>>> from urllib.request import urlopen
>>> data = np.loadtxt(urlopen("http://archive.ics.uci.edu/ml/machine-learning-databases/pima-indians-diabetes/
>>> pima-indians-diabetes.data"), delimiter=',')
>>> X = data[:, :-1]
>>> y = data[:, -1]
```

One-Hot Encoding

```
>>> from keras.utils import to_categorical
>>> y_train = to_categorical(y_train, num_classes)
>>> x_train = sequence.pad_sequences(x_train, maxlen=80)
>>> x_test4 = sequence.pad_sequences(x_test4, maxlen=80)
```

Preprocessing

Sequence Padding

```
>>> from keras.preprocessing import sequence
>>> x_train5, y_test5 = sequence.pad_sequences(x_train5, maxlen=80)
>>> x_test4 = sequence.pad_sequences(x_test4, maxlen=80)
```

Standardization/Normalization

```
>>> from keras.utils import to_categorical
>>> y_train3 = to_categorical(y_train3, num_classes)
>>> x_train3 = to_categorical(x_train3, num_classes)
>>> y_test3 = to_categorical(y_test3, num_classes)
```

Inspect Model

MLP: Binary Classification

```
>>> model.compile(optimizer='adam',
>>>                 loss='binary_crossentropy',
>>>                 metrics=['accuracy'])
```

MLP: Multi-Class Classification

```
>>> model.compile(optimizer='rmsprop',
>>>                 loss='categorical_crossentropy',
>>>                 metrics=['accuracy'])
```

MLP: Regression

```
>>> model.compile(optimizer='rmsprop',
>>>                 loss='mse',
>>>                 metrics=[mae])
```

Recurrent Neural Network

```
>>> model3.compile(loss='binary_crossentropy',
>>>                  optimizer='adam',
>>>                  metrics=['accuracy'])
```

Convolutional Neural Network (CNN)

```
>>> from keras.layers import Activation, Conv2D, MaxPooling2D, Flatten
>>> model2.add(Conv2D(32, (3, 3), padding='same', input_shape=x_train.shape[1:]))
>>> model2.add(Activation('relu'))
>>> model2.add(MaxPooling2D(pool_size=(2, 2)))
>>> model2.add(Dropout(0.25))
>>> model2.add(Flatten())
>>> model2.add(Dense(1))
```

Model Training

```
>>> model3.fit(x_train4,
>>>             y_train4,
>>>             batch_size=32,
>>>             epochs=5,
>>>             verbose=1,
>>>             validation_data=(x_test4, y_test4))
```

Evaluate Your Model's Performance

```
>>> score = model3.evaluate(x_test,
>>>                           y_test,
>>>                           batch_size=32)
```

Prediction

```
>>> model3.predict(x_test4, batch_size=32)
>>> model3.predict_classes(x_test4, batch_size=32)
```

Save/Reload Models

```
>>> from keras.models import load_model
>>> model3.save('model_file.h5')
>>> my_model = load_model('my_model.h5')
```

Model Fine-tuning

```
>>> model3.add(Embedding(20000, 128))
>>> model3.add(LSTM(128, dropout=0.2, recurrent_dropout=0.2))
>>> model3.add(Dense(1, activation='sigmoid'))
```

Optimization Parameters

```
>>> from keras.optimizers import RMSprop
>>> opt = RMSprop(lr=0.001, decay=1e-6)
>>> model2.compile(loss='categorical_crossentropy',
>>>                 optimizer=opt,
>>>                 metrics=['accuracy'])
```

Early Stopping

```
>>> from keras.callbacks import EarlyStopping
>>> early_stopping_monitor = EarlyStopping(patience=2)
>>> model3.fit(x_train4,
>>>             y_train4,
>>>             batch_size=32,
>>>             epochs=15,
>>>             validation_data=(x_test4, y_test4),
>>>             callbacks=[early_stopping_monitor])
```

Model output shape

```
Model summary
```

```
Model configuration
```

```
List all weight tensors in the model
```

TensorFlow v2.0 Cheat Sheet



TensorFlow™

TensorFlow is an open-source software library for high-performance numerical computation. Its flexible architecture enables to easily deploy computation across a variety of platforms (CPUs, GPUs, and TPUs), as well as mobile and edge devices, desktops, and clusters of servers. TensorFlow comes with strong support for machine learning and deep learning.

High-Level APIs for Deep Learning

Keras is a handy high-level API standard for deep learning models widely adopted for fast prototyping and state-of-the-art research. It was originally designed to run on top of different low-level computational frameworks and therefore the TensorFlow platform fully implements it.

The Sequential API is the most common way to define your neural network model. It corresponds to the mental image we use when thinking about deep learning: a sequence of layers.

```
import tensorflow as tf
from tensorflow.keras import datasets, layers, models

# Load data set
mnist = datasets.mnist
(x_train, y_train), (x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

# Construct a neural network model
model = models.Sequential()
model.add(layers.Flatten(input_shape=(28, 28)))
model.add(layers.Dense(512, activation=tf.nn.relu))
model.add(layers.Dropout(0.2))
model.add(layers.Dense(10, activation=tf.nn.softmax))
model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

# Train and evaluate the model
model.fit(x_train, y_train, epochs=5)
model.evaluate(x_test, y_test)
```

The Functional API enables engineers to define complex topologies, including multi-input and multi-output models, as well as advanced models with shared layers and models with residual connections.

```
from tensorflow.keras.layers import Flatten, Dense, Dropout
from tensorflow.keras.models import Model

# Loading data set must be here <...>
inputs = tf.keras.Input(shape=(28, 28))
x = Flatten()(inputs)
x = Dense(512, activation='relu')(x)
x = Dropout(0.2)(x)
predictions = Dense(10, activation='softmax')(x)
model = Model(inputs=inputs, outputs=predictions)

# Compile, train and evaluate the model here <...>
```

A layer instance is called on a tensor and returns a tensor. An input tensor and output tensor can then be used to define a Model, which is compiled and trained just as a Sequential model. Models are callable by themselves and can be stacked the same way while reusing trained weights.

Transfer learning and fine-tuning of pretrained models saves your time if your data set does not differ significantly from the original one.

```
import tensorflow as tf
import tensorflow_datasets as tfds

dataset = tfds.load(name='tf_flowers', as_supervised=True)
NUMBER_OF_CLASSES_IN_DATASET = 5
IMG_SIZE = 160

def preprocess_example(image, label):
    image = tf.cast(image, tf.float32)
    image = (image / 127.5) - 1
    image = tf.image.resize(image, (IMG_SIZE, IMG_SIZE))
    return image, label

DATASET_SIZE = 3670
BATCH_SIZE = 32
train = dataset['train'].map(preprocess_example)
train_batches = train.shuffle(DATASET_SIZE).batch(BATCH_SIZE)

# Load MobileNetV2 model pretrained on ImageNet data
model = tf.keras.applications.MobileNetV2(
    input_shape=(IMG_SIZE, IMG_SIZE, 3),
    include_top=False, weights='imagenet', pooling='avg')
model.trainable = False

# Add a new layer for multiclass classification
new_output = tf.keras.layers.Dense(
    NUMBER_OF_CLASSES_IN_DATASET, activation='softmax')
new_model = tf.keras.Sequential([model, new_output])
new_model.compile(
    loss=tf.keras.losses.categorical_crossentropy,
    optimizer=tf.keras.optimizers.RMSprop(lr=1e-3),
    metrics=['accuracy'])

# Train the classification layer
new_model.fit(train_batches.repeat(), epochs=10,
              steps_per_epoch=DATASET_SIZE // BATCH_SIZE)
```

After the execution of the given transfer learning code, you can make MobileNetV2 layers trainable and perform fine-tuning of the resulting model to achieve better results.

Jupyter Notebook

Jupyter Notebook is a web-based interactive computational environment for data science and scientific computing.

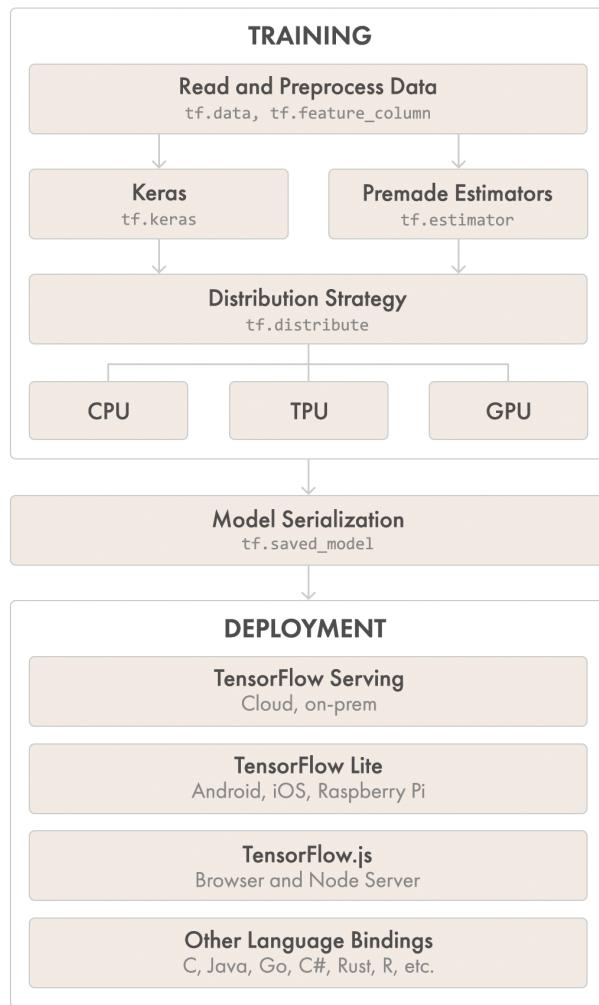
Google Colaboratory is a free notebook environment that requires no setup and runs entirely in the cloud. Use it for jump-starting a machine learning project.

TensorFlow v2.0 Cheat Sheet



A Reference Machine Learning Workflow

Here's a conceptual diagram and a workflow example:



01 Load the training data using **pipelines** created with `tf.data`. As an input, you can use either in-memory data (NumPy), or a local storage, or a remote persistent storage.

02 Build, train, and validate a model with `tf.keras`, or use premade estimators.

03 Run and debug with eager execution, then use `tf.function` for the benefits of graphs.

04 For large ML training tasks, use the **Distribution Strategy API** for deploying training on Kubernetes clusters within on-premises or cloud environments.

05 Export to **SavedModel**—an interchange format for TensorFlow Serving, TensorFlow Lite, TensorFlow.js, etc.

The `tf.data API` enables to build complex input pipelines from simple pieces. The pipeline aggregates data from a distributed file system, applies transformation to each object, and merges shuffled examples into training batches.

`tf.data.Dataset` represents a sequence of elements each containing one or more Tensor object(-s). This can be exemplified by a pair of tensors representing an image and a corresponding class label.

```
import tensorflow as tf

DATASET_URL = "https://archive.ics.uci.edu/ml/machine-learning-databases/covtype/covtype.data.gz"
DATASET_SIZE = 387698
dataset_path = tf.keras.utils.get_file(
    fname=DATASET_URL.split('/')[-1], origin=DATASET_URL)

COLUMN_NAMES = [
    'Elevation', 'Aspect', 'Slope',
    'Horizontal_Distance_To_Hydrology',
    'Vertical_Distance_To_Hydrology',
    'Horizontal_Distance_To_Roadways',
    'Hillshade_9am', 'Hillshade_Noon', 'Hillshade_3pm',
    'Horizontal_Distance_To_Fire_Points', 'Soil_Type',
    'Cover_Type']

def _parse_line(line):
    # Decode the line into values
    fields = tf.io.decode_csv(
        records=line, record_defaults=[0.0] * 54 + [0])

    # Pack the result into a dictionary
    features = dict(zip(COLUMN_NAMES,
        fields[:10] + [tf.stack(fields[14:54])] + [fields[-1]]))

    # Extract one-hot encoded class label from the features
    class_label = tf.argmax(features[10:14], axis=0)
    return features, class_label

def csv_input_fn(csv_path, test=False,
                 batch_size=DATASET_SIZE // 1000):
    # Create a dataset containing the csv lines
    dataset = tf.data.TextLineDataset(filenames=csv_path,
                                      compression_type='GZIP')

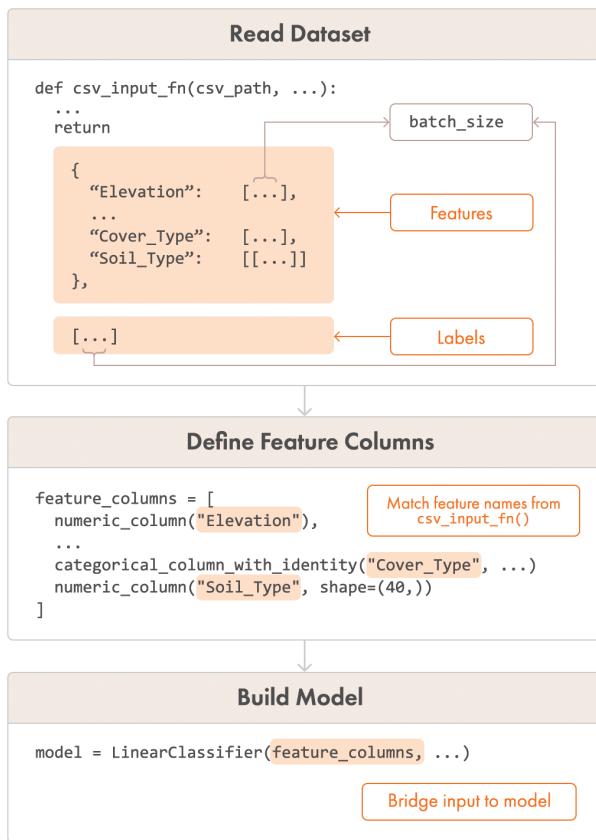
    # Parse each line
    dataset = dataset.map(_parse_line)

    # Shuffle, repeat, batch the examples for train and test
    dataset = dataset.shuffle(buffer_size=DATASET_SIZE,
                             seed=42)

    TEST_SIZE = DATASET_SIZE // 10
    return dataset.take(TEST_SIZE).batch(TEST_SIZE) if test \
        else dataset.skip(TEST_SIZE).repeat().batch(batch_size)
```

Functions from the `tf.feature_column` namespace are used to put raw data into a TensorFlow data set. A feature column is a high-level configuration abstraction for ingesting and representing features. It does not contain any data but tells the model how to transform the raw data so that it matches the expectation. The exact feature column to choose depends on the feature type and the model type. The continuous feature type is handled by `numeric_column` and can be directly fed into a neural network or a linear model.

TensorFlow v2.0 Cheat Sheet



Categorical features can be ingested by functions with the `"categorical_column"` prefix, but they need to be wrapped by `embedding_column` or `indicator_column` before being fed into Neural Network models. For linear models, `indicator_column` is an internal representation when categorical columns are passed in directly.

```
feature_columns = [tf.feature_column.numeric_column(name)
                   for name in COLUMN_NAMES[:10]]

feature_columns.append(
    tf.feature_column.categorical_column_with_identity(
        'Cover_Type', num_buckets=8)
)
# Soil_type[1-40] is a tensor of length 40
feature_columns.append(
    tf.feature_column.numeric_column('Soil_Type', shape=(40,))
)
```

The [Estimator API](#) provides high-level encapsulation for best practices: model training, evaluation, prediction, and export for serving. The `tf.estimator.Estimator` subclass represents a complete model. Its object creates and manages `tf.Graph` and `tf.Session` for you. [Premade estimators](#) include Linear Classifier, DNN Classifier, and Gradient Boosted Trees. `BaselineClassifier` and `BaselineRegressor` will help to establish a simple model for sanity check during further model development.

```
# Build, train, and evaluate the estimator
model = tf.estimator.LinearClassifier(feature_columns,
                                       n_classes=4)
model.train(input_fn=lambda: csv_input_fn(dataset_path),
            steps=10000)
model.evaluate(
    input_fn=lambda: csv_input_fn(dataset_path, test=True))
```

SavedModel contains a complete TF program and does not require the original model-building code to run, which makes it useful for deploying and sharing models.

```
# Export model to SavedModel
_builder = tf.estimator.export. \
    build_parsing_serving_input_receiver_fn
_spec_maker = tf.feature_column.make_parse_example_spec
serving_input_fn = _builder(_spec_maker(feature_columns))
export_path = model.export_saved_model(
    "/tmp/from_estimator/", serving_input_fn)
```

The following code sample shows how to load and use the saved model with Python.

```
# Import model from SavedModel
imported = tf.saved_model.load(export_path)

# Use imported model for prediction
def predict(new_object):
    example = tf.train.Example()

    # All regular continuous features
    for column in COLUMN_NAMES[:-2]:
        val = new_object[column]
        example.features.feature[column]. \
            float_list.value.extend([val])

    # One-hot encoded feature of 40 columns
    for val in new_object['Soil_Type']:
        example.features.feature['Soil_Type']. \
            float_list.value.extend([val])

    # Categorical column with ID
    example.features.feature['Cover_Type']. \
        int64_list.value.extend([new_object['Cover_Type']])

    return imported.signatures['predict'](
        examples=tf.constant([example.SerializeToString()]))

predict({
    'Elevation': 2296, 'Aspect': 312, 'Slope': 27,
    'Horizontal_Distance_To_Hydrology': 256,
    'Horizontal_Distance_To_Fire_Points': 836,
    'Horizontal_Distance_To_Roadways': 1273,
    'Vertical_Distance_To_Hydrology': 145,
    'Hillshade_9am': 136, 'Hillshade_Noon': 208,
    'Hillshade_3pm': 206,
    'Soil_Type': [0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
    'Cover_Type': 6})
```

VIP Cheatsheet: Tips and Tricks

Afishne AMIDI and Shervine AMIDI

November 26, 2018

Data processing



Data augmentation – Deep learning models usually need a lot of data to be properly trained.

It is often useful to get more data from the existing ones using data augmentation techniques. The main ones are summed up in the table below. More precisely, given the following input image, here are the techniques that we can apply:

| Original | Flip | Rotation | Random crop |
|----------------------------------|--|--|-------------|
| | | | |
| - Image without any modification | <ul style="list-style-type: none"> - Flipped with respect to an axis for which the meaning of the image is preserved - Rotation with a slight angle - Simulates incorrect horizon calibration | <ul style="list-style-type: none"> - Random focus on one part of the image - Several random crops can be done in a row | |

| Color shift | Noise addition | Information loss | Contrast change |
|-------------|----------------|------------------|-----------------|
| | | | |

Training a neural network



Epoch – In the context of training a model, epoch is a term used to refer to one iteration where the model sees the whole training set to update its weights.

Mini-batch gradient descent – During the training phase, updating weights is usually not based on the whole training set at once due to computation complexities or one data point due to noise issues. Instead, the update step is done on mini-batches, where the number of data points in a batch is a hyperparameter that we can tune.

Loss function – In order to quantify how a given model performs, the loss function L is usually used to evaluate to what extent the actual outputs y are correctly predicted by the model outputs z .

Cross-entropy loss – In the context of binary classification in neural networks, the cross-entropy loss $L(z,y)$ is commonly used and is defined as follows:

$$L(z,y) = - \left[y \log(z) + (1-y) \log(1-z) \right]$$

Backpropagation – Backpropagation is a method to update the weights in the neural network by taking into account the actual output and the desired output. The derivative with respect to each weight w is computed using the chain rule:

$$\frac{\partial L}{\partial f(x)} \cdot \frac{\partial f(x)}{\partial x} \cdot \frac{\partial x}{\partial f}$$

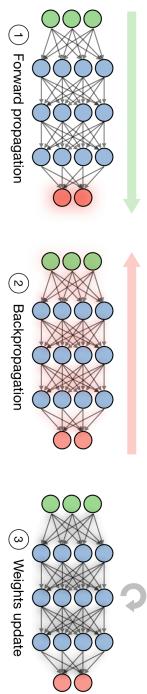
Using this method, each weight is updated with the rule:

$$w \leftarrow w - \alpha \frac{\partial L(z,y)}{\partial w}$$

Updating weights – In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data and perform forward propagation to compute the loss.
- Step 2: Backpropagate the loss to get the gradient of the loss with respect to each weight.
- Step 3: Use the gradients to update the weights of the network.

Batch normalization – It is a step of hyperparameter γ, β that normalizes the batch $\{x_i\}$. By noting μ_B, σ_B^2 the mean and variance of that we want to correct to the batch, it is done as follows:



Parameter tuning

Xavier initialization – Instead of initializing the weights in a purely random manner, Xavier initialization enables to have initial weights that take into account characteristics that are unique to the architecture.

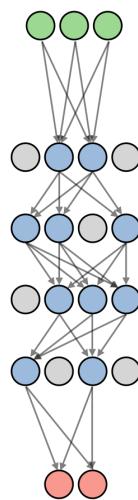
Transfer learning – Training a deep learning model requires a lot of data and more importantly a lot of time. It is often useful to take advantage of pre-trained weights on huge datasets that took days/weeks to train, and leverage it towards our use case. Depending on how much data we have at hand, here are the different ways to leverage this:

| Training size | Illustration | Explanation |
|---------------|--------------|--|
| Small | | Freezes all layers, trains weights on softmax |
| Medium | | Freezes most layers, trains weights on last layers and softmax |
| Large | | Trains weights on layers and softmax by initializing weights on pre-trained ones |

Remark: most deep learning frameworks parametrize dropout through the 'keep' parameter $1 - p$.

Dropout – Dropout is a technique used in neural networks to prevent overfitting the training data by dropping out neurons with probability $p > 0$. It forces the model to avoid relying too much on particular sets of features.

Regularization



Remark: other methods include Adadelta, Adagrad and SGD.

| Method | Explanation | Update of w | Update of b |
|----------|---|---|---|
| Momentum | - Dampens oscillations - Improvement to SGD - 2 parameters to tune | $w \leftarrow w - \alpha v_{dw}$ | $b \leftarrow b - \alpha v_{db}$ |
| RMSprop | - Root Mean Square propagation by controlling oscillations | $w \leftarrow w - \alpha \frac{dw}{\sqrt{s_{dw}}}$ | $b \leftarrow b - \alpha \frac{db}{\sqrt{s_{db}}}$ |
| Adam | - Adaptive Moment estimation - Most popular method - 4 parameters to tune | $w \leftarrow w - \alpha \frac{v_{dw}}{\sqrt{s_{dw}} + \epsilon}$ | $b \leftarrow b - \alpha \frac{v_{db}}{\sqrt{s_{db}} + \epsilon}$ |

- Learning rate** – The learning rate, often noted α or sometimes η , indicates at which pace the weights get updated. It can be fixed or adaptively changed. The current most popular method is called Adam, which is a method that adapts the learning rate.
- Adaptive learning rates** – Letting the learning rate vary when training a model can reduce the training time and improve the numerical optimal solution. While Adam optimizer is the most commonly used technique, others can also be useful. They are summed up in the table below:

| LASSO | Ridge | Elastic Net |
|--|----------------------------|--|
| <ul style="list-style-type: none"> - Shrinks coefficients to 0 - Good for variable selection | Makes coefficients smaller | Tradeoff between variable selection and small coefficients |

$$\dots + \lambda \|\theta\|_1$$

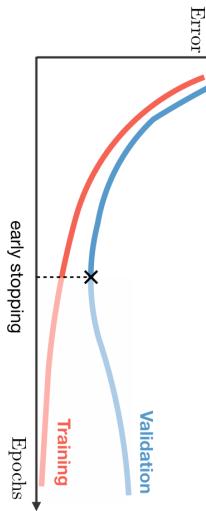
$$\lambda \in \mathbb{R}$$

$$(1 - \alpha)\|\theta\|_1 + \alpha\|\theta\|_2^2 \leq 1$$

$$\dots + \lambda \|\theta\|_2^2$$

$$\lambda \in \mathbb{R}, \alpha \in [0,1]$$

- **Early stopping** – This regularization technique stops the training process as soon as the validation loss reaches a plateau or starts to increase.



Good practices

- **Overfitting small batch** – When debugging a model, it is often useful to make quick tests to see if there is any major issue with the architecture of the model itself. In particular, in order to make sure that the model can be properly trained, a mini-batch is passed inside the network to see if it can overfit on it. If it cannot, it means that the model is either too complex or not complex enough to even overfit on a small batch, let alone a normal-sized training set.
- **Gradient checking** – Gradient checking is a method used during the implementation of the backward pass of a neural network. It compares the value of the analytical gradient to the numerical gradient at given points and plays the role of a sanity-check for correctness.

| | Numerical gradient | Analytical gradient |
|----------|--|--|
| Formula | $\frac{df}{dx}(x) \approx \frac{f(x + h) - f(x - h)}{2h}$ | $\frac{df}{dx}(x) = f'(x)$ |
| Comments | <ul style="list-style-type: none"> - Expensive; loss has to be computed two times per dimension - Used to verify correctness of analytical implementation - Trade-off in choosing h: not too small (numerical instability) nor too large (poor gradient approx.) | <ul style="list-style-type: none"> - 'Exact' result - Direct computation - Used in the final implementation |

* * *

VIP Cheatsheet: Deep Learning

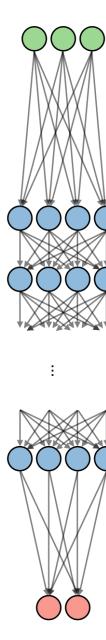
Afshine AMIDI and Shervine AMIDI

September 15, 2018

Neural Networks

Neural networks are a class of models that are built with layers. Commonly used types of neural networks include convolutional and recurrent neural networks.

Architecture – The vocabulary around neural networks architectures is described in the figure below:



By noting i the i^{th} layer of the network and j the j^{th} hidden unit of the layer, we have:

$$z_j^{[i]} = w_j^{[i]T} x + b_j^{[i]}$$

where we note w , b , z the weight, bias and output respectively.

Activation function – Activation functions are used at the end of a hidden unit to introduce non-linear complexities to the model. Here are the most common ones:

| Sigmoid | Tanh | ReLU | Leaky ReLU |
|-------------------------------|--|---------------------|---|
| $g(z) = \frac{1}{1 + e^{-z}}$ | $g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ | $g(z) = \max(0, z)$ | $g(z) = \max(\epsilon z, z)$ with $\epsilon \ll 1$ |

Convolutional Neural Networks

Convolutional layer requirement – By noting W the input volume size, F the size of the convolutional layer neurons, P the amount of zero padding, then the number of neurons N that fit in a given volume is such that:

$$N = \frac{W - F + 2P}{S} + 1$$

Batch normalization – It is a step of hyperparameter γ, β that normalizes the batch $\{x_i\}$. By noting μ_B, σ_B^2 the mean and variance of that we want to correct to the batch, it is done as follows:

$$x_i \leftarrow \gamma \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

Cross-entropy loss – In the context of neural networks, the cross-entropy loss $L(z,y)$ is commonly used and is defined as follows:

$$L(z,y) = - \left[y \log(z) + (1-y) \log(1-z) \right]$$

Learning rate – The learning rate, often noted η ; indicates at which pace the weights get updated. This can be fixed or adaptively changed. The current most popular method is called Adam, which is a method that adapts the learning rate.

Backpropagation – Backpropagation is a method to update the weights in the neural network by taking into account the actual output and the desired output. The derivative with respect to weight w is computed using chain rule and is of the following form:

$$\frac{\partial L(z,y)}{\partial w} = \frac{\partial L(z,y)}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w}$$

As a result, the weight is updated as follows:

$$w \leftarrow w - \eta \frac{\partial L(z,y)}{\partial w}$$

Updating weights – In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data.
- Step 2: Perform forward propagation to obtain the corresponding loss.
- Step 3: Backpropagate the loss to get the gradients.
- Step 4: Use the gradients to update the weights of the network.

It is usually done after a fully connected/convolutional layer and before a non-linearity layer and aims at allowing higher learning rates and reducing the strong dependence on initialization.

Recurrent Neural Networks

□ Types of gates – Here are the different types of gates that we encounter in a typical recurrent neural network:

| Input gate | Forget gate | Output gate | Gate |
|-----------------------|----------------------|-----------------------|-------------------|
| Write to cell or not? | Erase a cell or not? | Reveal a cell or not? | How much writing? |

□ LSTM – A long short-term memory (LSTM) network is a type of RNN model that avoids the vanishing gradient problem by adding 'forget' gates.

Reinforcement Learning and Control

The goal of reinforcement learning is for an agent to learn how to evolve in an environment.

□ Markov decision processes – A Markov decision process (MDP) is a 5-tuple $(S, A, \{P_{sa}\}, \gamma, R)$ where:

- S is the set of states
- A is the set of actions
- $\{P_{sa}\}$ are the state transition probabilities for $s \in S$ and $a \in A$
- $\gamma \in [0,1]$ is the discount factor
- $R : S \times A \rightarrow \mathbb{R}$ or $R : S \rightarrow \mathbb{R}$ is the reward function that the algorithm wants to maximize

□ Policy – A policy π is a function $\pi : S \rightarrow A$ that maps states to actions.

Remark: we say that we execute a given policy π if given a state s we take the action $a = \pi(s)$.

□ Value function – For a given policy π and a given state s , we define the value function V^π as follows:

$$V^\pi(s) = E \left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | s_0 = s, \pi \right]$$

□ Bellman equation – The optimal Bellman equations characterizes the value function V^{π^*} of the optimal policy π^* :

$$V^{\pi^*}(s) = R(s) + \max_{a \in A} \gamma \sum_{s' \in S} P_{sa}(s') V^{\pi^*}(s')$$

Remark: we note that the optimal policy π^ for a given state s is such that:*

$$\pi^*(s) = \arg\max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

□ Value iteration algorithm – The value iteration algorithm is in two steps:

- We initialize the value:

$$V_0(s) = 0$$

- We iterate the value based on the values before:

$$V_{i+1}(s) = R(s) + \max_{a \in A} \left[\sum_{s' \in S} \gamma P_{sa}(s') V_i(s') \right]$$

□ Maximum likelihood estimate – The maximum likelihood estimates for the state transition probabilities are as follows:

$$P_{sa}(s') = \frac{\# \text{times took action } a \text{ in state } s \text{ and got to } s'}{\# \text{times took action } a \text{ in state } s}$$

□ Q-learning – Q-learning is a model-free estimation of Q , which is done as follows:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

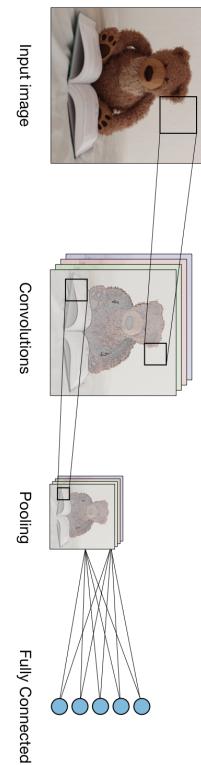
VIP Cheatsheet: Convolutional Neural Networks

Afshine AMIDI and Shervine AMIDI

November 26, 2018

Overview

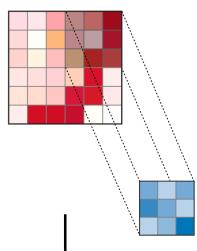
□ **Architecture of a traditional CNN** – Convolutional neural networks, also known as CNNs, are a specific type of neural networks that are generally composed of the following layers:



The convolution layer and the pooling layer can be fine-tuned with respect to hyperparameters that are described in the next sections.

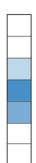
Types of layer

□ **Convolutional layer (CONV)** – The convolution layer (CONV) uses filters that perform convolution operations as it is scanning the input I with respect to its dimensions. Its hyperparameters include the filter size F and stride S . The resulting output O is called *feature map* or *activation map*.



Remark: the convolution step can be generalized to the 1D and 3D cases as well.

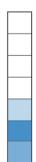
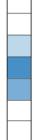
□ **Pooling (POOL)** – The pooling layer (POOL) is a downsampling operation, typically applied after a convolution layer, which does some spatial invariance. In particular, max and average pooling are special kinds of pooling where the maximum and average value is taken, respectively.



Remark: the application of K filters of size $F \times F$ results in an output feature map of size $O \times O \times K$.

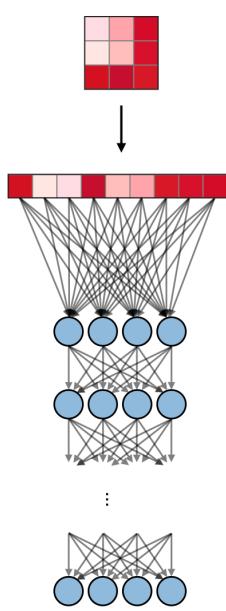
Remark: the convolution step can be generalized to the 1D and 3D cases as well.

□ **Stride** – For a convolutional or a pooling operation, the stride S denotes the number of pixels by which the window moves after each operation.



| Purpose | Max pooling | Average pooling |
|--------------|---|--|
| | Each pooling operation selects the maximum value of the current view | Each pooling operation averages the values of the current view |
| Illustration | | |
| Comments | <ul style="list-style-type: none"> - Preserves detected features - Most commonly used | <ul style="list-style-type: none"> - Downsamples feature map - Used in LeNet |

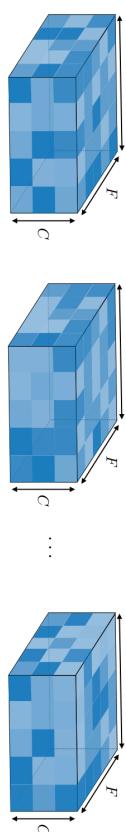
□ **Fully Connected (FC)** – The fully connected layer (FC) operates on a flattened input where each input is connected to all neurons. If present, FC layers are usually found towards the end of CNN architectures and can be used to optimize objectives such as class scores.



Filter hyperparameters

The convolution layer contains filters for which it is important to know the meaning behind its hyperparameters.

□ **Dimensions of a filter** – A filter of size $F \times F$ applied to an input I containing C channels is a $F \times F \times C$ volume that performs convolutions on an input of size $I \times I \times C$ and produces an output feature map (also called activation map) of size $O \times O \times 1$.



Remark: the application of K filters of size $F \times F$ results in an output feature map of size $O \times O \times K$.

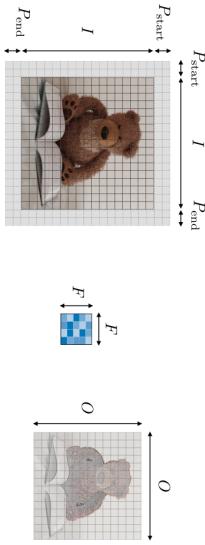
□ **Zero-padding** – Zero-padding denotes the process of adding P zeroes to each side of the boundaries of the input. This value can either be manually specified or automatically set through one of the three modes detailed below:

| | Valid | Same | Full |
|---------------------|---|---|---|
| Value | $P = 0$ | $P_{\text{start}} = \left\lceil \frac{S\lceil \frac{I}{S} \rceil - I + F - S}{2} \right\rceil$ $P_{\text{end}} = \left\lceil \frac{S\lceil \frac{I}{S} \rceil - I + F - S}{2} \right\rceil$ | $P_{\text{start}} \in [0, F - 1]$ $P_{\text{end}} = F - 1$ |
| Illustration | | | |
| Purpose | <ul style="list-style-type: none"> - No padding - Drops last convolution if dimensions do not match | <ul style="list-style-type: none"> - Padding such that feature map size has size $\left\lceil \frac{I}{S} \right\rceil$ - Output size is mathematically convenient - Also called 'half' padding | <ul style="list-style-type: none"> - Maximum padding such that end convolutions are applied on the limits of the input - Filter 'sees' the input end-to-end |

Tuning hyperparameters

□ **Parameter compatibility** in convolution layer – By noting I the length of the input volume size, F the length of the filter, P the amount of zero padding, S the stride, then the output size O of the feature map along that dimension is given by:

$$O = \frac{I - F + P_{\text{start}} + P_{\text{end}}}{S} + 1$$



Remark: often times, $P_{\text{start}} = P_{\text{end}} \triangleq P$, in which case we can replace $P_{\text{start}} + P_{\text{end}}$ by $2P$ in the formula above.

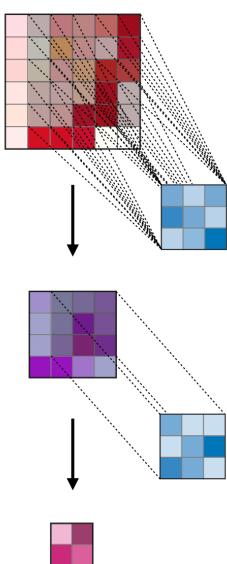
□ **Understanding the complexity of the model** – In order to assess the complexity of a model, it is often useful to determine the number of parameters that its architecture will have. In a given layer of a convolutional neural network, it is done as follows:

| | CONV | POOL | FC |
|-----------------------------|--|--|---|
| Illustration | | | |
| Input size | $I \times I \times C$ | $I \times I \times C$ | N_{in} |
| Output size | $O \times O \times K$ | $O \times O \times C$ | N_{out} |
| Number of parameters | $(F \times F \times C + 1) \cdot K$ | 0 | $(N_{\text{in}} + 1) \times N_{\text{out}}$ |
| Remarks | <ul style="list-style-type: none"> - One bias parameter per filter - In most cases, $S < F$ - A common choice for K is $2C$ | <ul style="list-style-type: none"> - Pooling operation done channel-wise - In most cases, $S = F$ | <ul style="list-style-type: none"> - Input is flattened - One bias parameter per neuron - The number of FC neurons is free of structural constraints |

□ **Receptive field** – The receptive field at layer k is the area denoted $R_k \times R_k$ of the input that each pixel of the k -th activation map can see. By calling F_j the filter size of layer j and S_i the stride value of layer i and with the convention $S_0 = 1$, the receptive field at layer k can be computed with the formula:

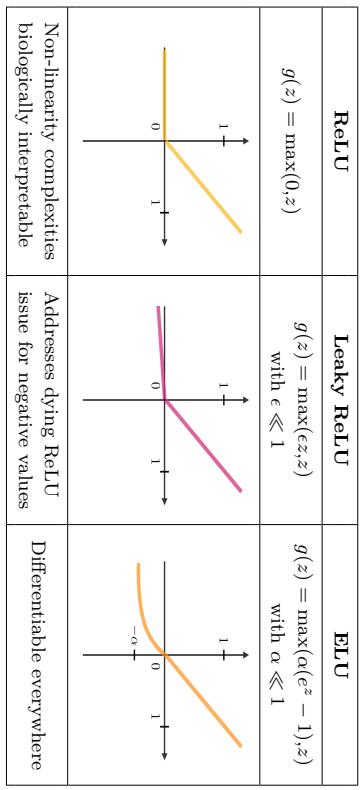
$$R_k = 1 + \sum_{j=1}^k (F_j - 1) \prod_{i=0}^{j-1} S_i$$

In the example below, we have $F_1 = F_2 = 3$ and $S_1 = S_2 = 1$, which gives $R_2 = 1+2 \cdot 1+2 \cdot 1 = 5$.



Commonly used activation functions

□ **Rectified Linear Unit** – The rectified linear unit layer (ReLU) is an activation function g that is used on all elements of the volume. It aims at introducing non-linearities to the network. Its variants are summarized in the table below:



□ **Softmax** – The softmax step can be seen as a generalized logistic function that takes as input a vector of scores $x \in \mathbb{R}^n$ and outputs a vector of output probability $p \in \mathbb{R}^n$ through a softmax function at the end of the architecture. It is defined as follows:

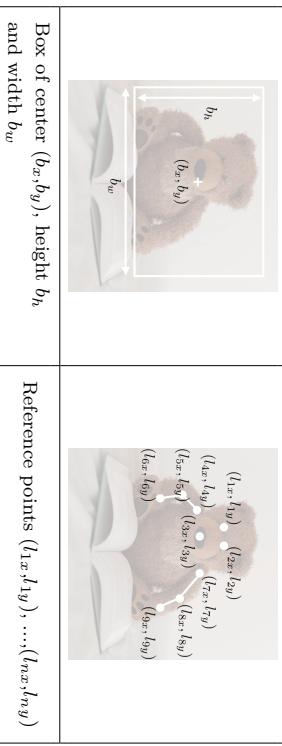
$$p = \begin{pmatrix} p_1 \\ \vdots \\ p_n \end{pmatrix} \quad \text{where} \quad p_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

Object detection

□ **Types of models** – There are 3 main types of object recognition algorithms, for which the nature of what is predicted is different. They are described in the table below:

| Image classification | Classification w. localization | Detection |
|----------------------|--------------------------------|-------------|
| | | |
| | | |
| Traditional CNN | Simplified YOLO, R-CNN | YOLO, R-CNN |

| Bounding box detection | Landmark detection |
|---|--|
| Detects the part of the image where the object is located | - Detects a shape or characteristics of an object (e.g. eyes) - More granular |



□ **Intersection over Union** – Intersection over Union, also known as IoU, is a function that quantifies how correctly positioned a predicted bounding box B_p is over the actual bounding box B_a . It is defined as:

$$\text{IoU}(B_p, B_a) = \frac{B_p \cap B_a}{B_p \cup B_a}$$



Remark: we always have $\text{IoU} \in [0,1]$. By convention, a predicted bounding box B_p is considered as being reasonably good if $\text{IoU}(B_p, B_a) \geq 0.5$.

□ **Anchor boxes** – Anchor boxing is a technique used to predict overlapping bounding boxes. In practice, the network is allowed to predict more than one box simultaneously, where each box prediction is constrained to have a given set of geometrical properties. For instance, the first prediction can potentially be a rectangular box of a given form, while the second will be another rectangular box of a different geometrical form.

□ **Non-max suppression** – The non-max suppression technique aims at removing duplicate overlapping bounding boxes of a same object by selecting the most representative ones. After having removed all boxes having a probability prediction lower than 0.6, the following steps are repeated while there are boxes remaining:

- Step 1: Pick the box with the largest prediction probability.
- Step 2: Discard any box having an $\text{IoU} \geq 0.5$ with the previous box.



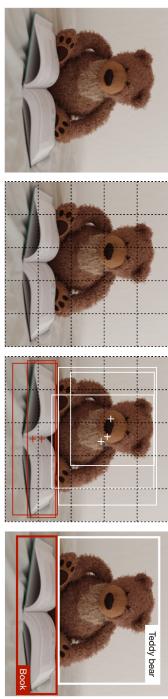
□ YOLO – You Only Look Once (YOLO) is an object detection algorithm that performs the following steps:

- Step 1: Divide the input image into a $G \times G$ grid.
- Step 2: For each grid cell, run a CNN that predicts y of the following form:

$$y = \underbrace{\begin{bmatrix} p_c, b_x, b_y, b_h, b_w, c_1, c_2, \dots, c_p, \dots \end{bmatrix}}_{\text{repeated } k \text{ times}}^T \in \mathbb{R}^{G \times G \times k \times (5+p)}$$

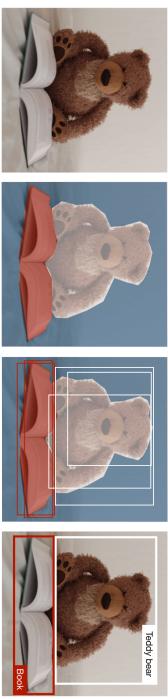
where p_c is the probability of detecting an object, b_x, b_y, b_h, b_w are the properties of the detected bounding box, c_1, \dots, c_p is a one-hot representation of which of the p classes were detected, and k is the number of anchor boxes.

- Step 3: Run the non-max suppression algorithm to remove any potential duplicate overlapping bounding boxes.



Remark: when $p_c = 0$, then the network does not detect any object. In that case, the corresponding predictions b_x, \dots, c_p have to be ignored.

□ R-CNN – Region with Convolutional Neural Networks (R-CNN) is an object detection algorithm that first segments the image to find potential relevant bounding boxes and then run the detection algorithm to find most probable objects in those bounding boxes.



Face verification and recognition

□ Types of models – Two main types of model are summed up in table below:

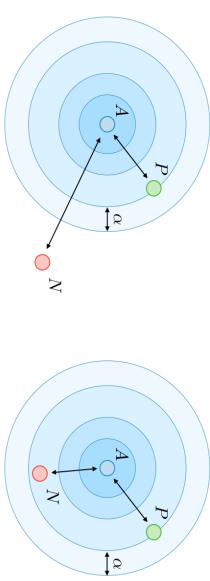
| Face verification | Face recognition |
|--|---|
| - Is this the correct person? - One-to-one lookup | - Is this one of the K persons in the database? - One-to-many lookup |
| | |

□ One Shot Learning – One Shot Learning is a face verification algorithm that uses a limited training set to learn a similarity function that quantifies how different two given images are. The similarity function applied to two images is often noted $d(\text{image 1}, \text{image 2})$.

□ Siamese Network – Siamese Networks aim at learning how to encode images to then quantify how different two images are. For a given input image $x^{(i)}$, the encoded output is often noted as $f(x^{(i)})$.

□ Triplet loss – The triplet loss ℓ is a loss function computed on the embedding representation of a triplet of images A (anchor), P (positive) and N (negative). The anchor and the positive example belong to a same class, while the negative example to another one. By calling $\alpha \in \mathbb{R}^+$ the margin parameter, this loss is defined as follows:

$$\ell(A, P, N) = \max(d(A, P) - d(A, N) + \alpha, 0)$$



Neural style transfer

Remark: although the original algorithm is computationally expensive and slow, newer architectures enabled the algorithm to run faster, such as Fast R-CNN and Faster R-CNN.

Original image

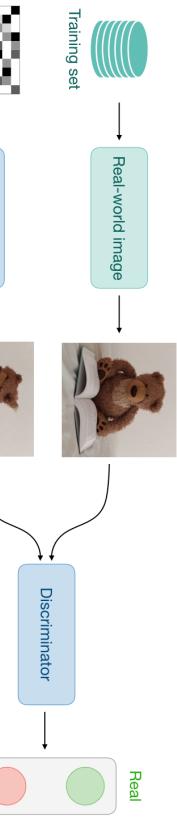
Segmentation

Bounding box prediction

Non-max suppression

$$\text{Content } C + \text{Style } S = \text{Generated image } G$$

The diagram illustrates the process of image generation. It shows three images: a teddy bear (Content C), a starry night painting (Style S), and a generated image G. A plus sign (+) is placed between Content C and Style S, followed by an equals sign (=) and the Generated image G.



- **Activation** – In a given layer l , the activation is noted $a^{[l]}$ and is of dimensions $n_H \times n_w \times n_c$
- **Content cost function** – The content cost function $J_{\text{content}}(C, G)$ is used to determine how the generated image G differs from the original content image C . It is defined as follows:

$$J_{\text{content}}(C, G) = \frac{1}{2} \|a^{[l](C)} - a^{[l](G)}\|^2$$

- **Style matrix** – The style matrix $G^{[l]}$ of a given layer l is a Gram matrix where each of its elements $G_{k k'}^{[l]}$ quantifies how correlated the channels k and k' are. It is defined with respect to activations $a^{[l]}$ as follows:

$$G_{k k'}^{[l]} = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} a_{ijk'}^{[l]}$$

Remark: the style matrix for the style image and the generated image are noted $G^{[l](S)}$ and $G^{[l](G)}$ respectively.

- **Style cost function** – The style cost function $J_{\text{style}}(S, G)$ is used to determine how the generated image G differs from the style S . It is defined as follows:

$$J_{\text{style}}^{[l]}(S, G) = \frac{1}{(2n_H n_w n_c)^2} \|G^{[l](S)} - G^{[l](G)}\|_F^2 = \frac{1}{(2n_H n_w n_c)^2} \sum_{k, k'=1}^{n_c} \left(G_{k k'}^{[l](S)} - G_{k k'}^{[l](G)} \right)^2$$

- **Overall cost function** – The overall cost function is defined as being a combination of the content and style cost functions, weighted by parameters α, β , as follows:

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

Remark: a higher value of α will make the model care more about the content while a higher value of β will make it care more about the style.

Architectures using computational tricks

- **Generative Adversarial Network** – Generative adversarial networks, also known as GANs, are composed of a generative and a discriminative model, where the generative model aims at generating the most truthful output that will be fed into the discriminative which aims at differentiating the generated and true image.

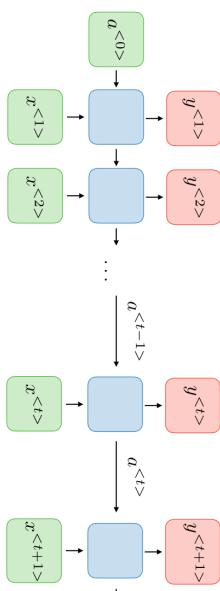
VIP Cheatsheet: Recurrent Neural Networks

Afshine AMIDI and Shervine AMIDI

November 26, 2018

Overview

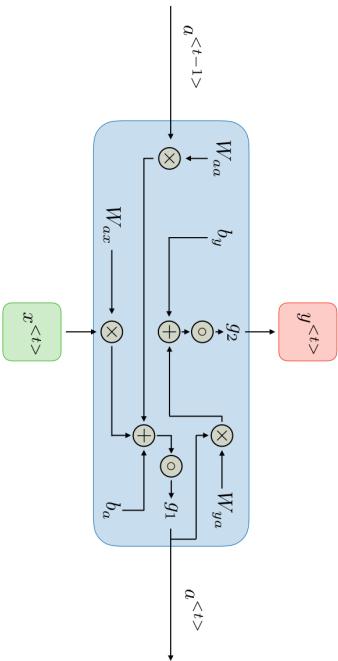
□ Architecture of a traditional RNN – Recurrent neural networks, also known as RNNs, are a class of neural networks that allow previous outputs to be used as inputs while having hidden states. They are typically as follows:



For each timestep t , the activation $a^{<t>}$ and the output $y^{<t>}$ are expressed as follows:

$$a^{<t>} = g_1(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a) \quad \text{and} \quad y^{<t>} = g_2(W_{ya}a^{<t>} + b_y)$$

where $W_{ax}, W_{aa}, W_{ya}, b_a, b_y$ are coefficients that are shared temporally and g_1, g_2 activation functions



| Type of RNN | Illustration | Example |
|-----------------------------------|--------------|----------------------------|
| One-to-one $T_x = T_y = 1$ | | Traditional neural network |
| One-to-many $T_x = 1, T_y > 1$ | | Music generation |
| Many-to-one $T_x > 1, T_y = 1$ | | Sentiment classification |
| Many-to-many $T_x = T_y$ | | Name entity recognition |
| Many-to-many $T_x \neq T_y$ | | Machine translation |

□ Applications of RNNs – RNN models are mostly used in the fields of natural language processing and speech recognition. The different applications are summed up in the table below:

The pros and cons of a typical RNN architecture are summed up in the table below:

steps is defined based on the loss at every time step as follows:

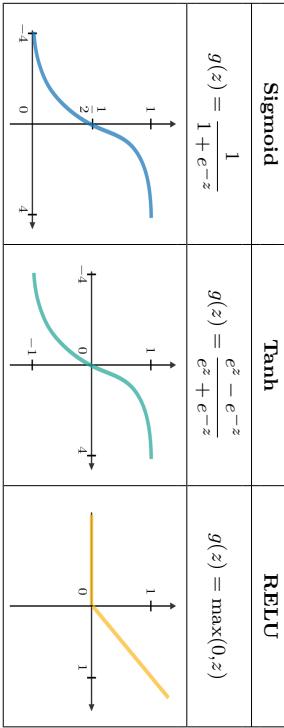
$$\mathcal{L}(\hat{y}) = \sum_{t=1}^{T_y} \mathcal{L}(\hat{y}^{<t>}, y^{<t>})$$

□ **Backpropagation through time** – Backpropagation is done at each point in time. At timestep T , the derivative of the loss \mathcal{L} with respect to weight matrix W is expressed as follows:

$$\frac{\partial \mathcal{L}^{(T)}}{\partial W} = \sum_{t=1}^T \left. \frac{\partial \mathcal{L}^{(T)}}{\partial W} \right|_{(t)}$$

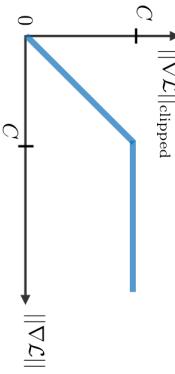
Handling long term dependencies

□ **Commonly used activation functions** – The most common activation functions used in RNN modules are described below:



□ **Vanishing/exploding gradient** – The vanishing and exploding gradient phenomena are often encountered in the context of RNNs. The reason why they happen is that it is difficult to capture long term dependencies because of multiplicative gradient that can be exponentially decreasing/increasing with respect to the number of layers.

□ **Gradient clipping** – It is a technique used to cope with the exploding gradient problem sometimes encountered when performing backpropagation. By capping the maximum value for the gradient, this phenomenon is controlled in practice.



Remark: the sign $$ denotes the element-wise multiplication between two vectors.*

□ **Variants of RNNs** – The table below sums up the other commonly used RNN architectures:

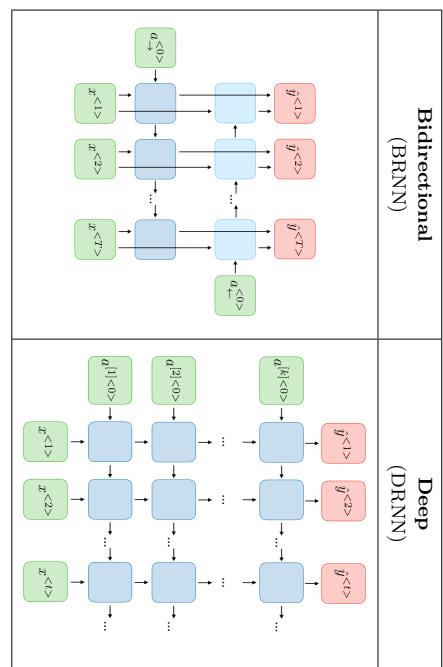
□ **Types of gates** – In order to remedy the vanishing gradient problem, specific gates are used in some types of RNNs and usually have a well-defined purpose. They are usually noted Γ and are equal to:

$$\Gamma = \sigma(Wx^{<t>} + Ua^{<t-1>} + b)$$

| Type of gate | Role | Used in |
|---------------------------|----------------------------------|-----------|
| Update gate Γ_u | How much past should matter now? | GRU, LSTM |
| Relevance gate Γ_r | Drop previous information? | GRU, LSTM |
| Forget gate Γ_f | Erase a cell or not? | LSTM |
| Output gate Γ_o | How much to reveal of a cell? | LSTM |

□ **GRU/LSTM** – Gated Recurrent Unit (GRU) and Long Short-Term Memory units (LSTM) deal with the vanishing gradient problem encountered by traditional RNNs, with LSTM being a generalization of GRU. Below is a table summing up the characterizing equations of each architecture:

| Gated Recurrent Unit (GRU) | | Long Short-Term Memory (LSTM) | |
|----------------------------|---|---|--|
| $\tilde{c}^{<t>}$ | $\tanh(W_c[\Gamma_r * a^{<t-1>}, x^{<t>}] + b_c)$ | $\tanh(W_c[\Gamma_r * a^{<t-1>}, x^{<t>}] + b_c)$ | |
| $c^{<t>}$ | $\Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$ | $\Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$ | |
| $a^{<t>}$ | $c^{<t>}$ | $\Gamma_o * c^{<t>}$ | |
| Dependencies | | | |



Learning word representation

In this section, we note V the vocabulary and $|V|$ its size.

□ Representation techniques – The two main ways of representing words are summed up in the table below:

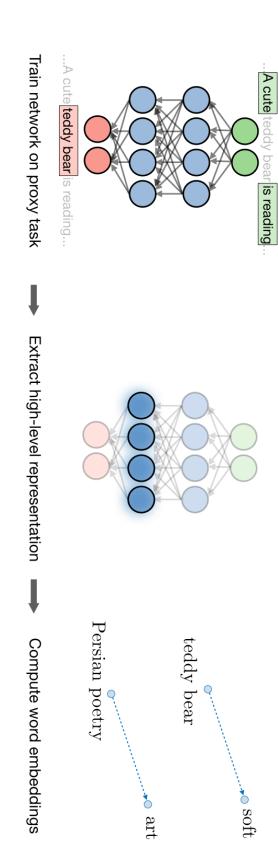
| 1-hot representation | Word embedding |
|-----------------------------------|---|
| <pre> teddy bear book soft </pre> | <pre> teddy bear teddy bear teddy bear </pre> |

- Noted o_w
- Naive approach, no similarity information

Remark: summing over the whole vocabulary in the denominator of the softmax part makes this model computationally expensive. CBOW is another word2vec model using the surrounding words to predict a given word.

□ Skip-gram – The skip-gram word2vec model is a supervised learning task that learns word embeddings by assessing the likelihood of any given target word t happening with a context word c . By noting θ_t a parameter associated with t , the probability $P(t|c)$ is given by:

$$P(t|c) = \frac{\exp(\theta_t^T e_c)}{\sum_{j=1}^{|V|} \exp(\theta_j^T e_c)}$$



Remark: this method is less computationally expensive than the skip-gram model.

□ GloVe – The GloVe model, short for global vectors for word representation, is a word embedding technique that uses a co-occurrence matrix X where each $X_{i,j}$ denotes the number of times that a target i occurred with a context j . Its cost function J is as follows:

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^{|V|} f(X_{i,j})(\theta_i^T e_j + b_i + b'_j - \log(X_{i,j}))^2$$

here f is a weighting function such that $X_{i,j} = 0 \implies f(X_{i,j}) = 0$. Given the symmetry that e and θ play in this model, the final word embedding $e_w^{(\text{final})}$ is given by:

$$e_w^{(\text{final})} = \frac{e_w + \theta_w}{2}$$

Remark: learning the embedding matrix can be done using target/context likelihood models.

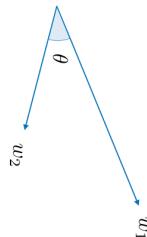
□ Word2vec – Word2vec is a framework aimed at learning word embeddings by estimating the likelihood that a given word is surrounded by other words. Popular models include skip-gram, negative sampling and CBOW.

Comparing words

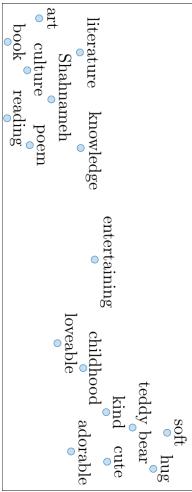
□ **Cosine similarity** – The cosine similarity between words w_1 and w_2 is expressed as follows:

$$\text{similarity} = \frac{w_1 \cdot w_2}{\|w_1\| \|w_2\|} = \cos(\theta)$$

Remark: θ is the angle between words w_1 and w_2 .



□ **t-SNE** – t-SNE (t -distributed Stochastic Neighbor Embedding) is a technique aimed at reducing high-dimensional embeddings into a lower dimensional space. In practice, it is commonly used to visualize word vectors in the 2D space.



Language model

□ **Overview** – A language model aims at estimating the probability of a sentence $P(y)$.

□ **n-gram model** – This model is a naive approach aiming at quantifying the probability that an expression appears in a corpus by counting its number of appearance in the training data.

□ **Perplexity** – Language models are commonly assessed using the perplexity metric, also known as PP, which can be interpreted as the inverse probability of the dataset normalized by the number of words T . The perplexity is such that the lower, the better and is defined as follows:

$$\text{PP} = \prod_{t=1}^T \left(\frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \hat{y}_j^{(t)}} \right)^{\frac{1}{T}}$$

Remark: PP is commonly used in t-SNE.

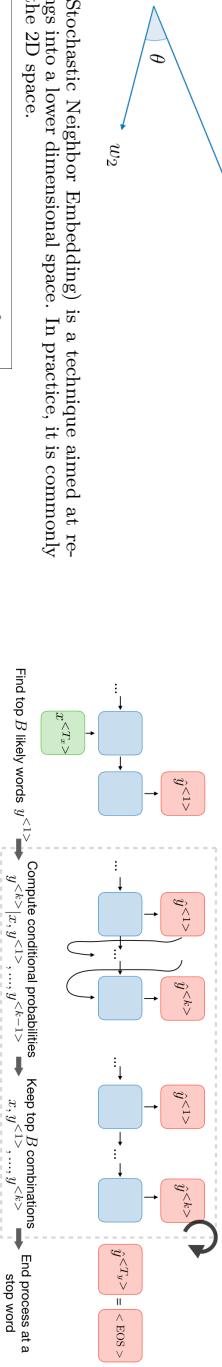
Machine translation

□ **Overview** – A machine translation model is similar to a language model except it has an encoder network placed before. For this reason, it is sometimes referred as a conditional language model. The goal is to find a sentence y such that:

$$y = \arg \max_{y^{<1>} \dots, y^{<T_y>} | x} P(y^{<1>} \dots, y^{<T_y>} | x)$$

□ **Beam search** – It is a heuristic search algorithm used in machine translation and speech recognition to find the likelihood sentence y given an input x .

- Step 1: Find top B likely words $y^{<1>}$
- Step 2: Compute conditional probabilities $y^{<k>} | x, y^{<1>} \dots, y^{<k-1>}$
- Step 3: Keep top B combinations $x, y^{<1>} \dots, y^{<k>}$



Remark: if the beam width is set to 1, then this is equivalent to a naive greedy search.

□ **Beam width** – The beam width B is a parameter for beam search. Large values of B yield to better result but with slower performance and increased memory. Small values of B lead to worse results but is less computationally intensive. A standard value for B is around 10.

□ **Length normalization** – In order to improve numerical stability, beam search is usually applied on the following normalized objective, often called the normalized log-likelihood objective, defined as:

$$\text{Objective} = \frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log \left[p(y^{<t>} | x, y^{<1>} \dots, y^{<t-1>}) \right]$$

Remark: the parameter α can be seen as a softener, and its value is usually between 0.5 and 1.

□ **Error analysis** – When obtaining a predicted translation \hat{y} that is bad, one can wonder why we did not get a good translation y^* by performing the following error analysis:

| Case | $P(y^* x) > P(\hat{y} x)$ | $P(y^* x) \leqslant P(\hat{y} x)$ |
|------------|--|---------------------------------------|
| Root cause | Beam search faulty | RNN faulty |
| Remedies | <ul style="list-style-type: none"> - Increase beam width - Regularize - Get more data | |

□ **Bleu score** – The bilingual evaluation under study (bleu) score quantifies how good a machine translation is by computing a similarity score based on n -gram precision. It is defined as follows:

$$\text{bleu score} = \exp \left(\frac{1}{n} \sum_{k=1}^n p_k \right)$$

where p_n is the bleu score on n -gram only defined as follows:

$$p_n = \frac{\sum_{\text{n-gram} \in \hat{y}} \text{count}_{\text{clip}}(\text{n-gram})}{\sum_{\text{n-gram} \in \hat{y}} \text{count}(\text{n-gram})}$$

Remark: a brevity penalty may be applied to short predicted translations to prevent an artificially inflated bleu score.

Attention

□ **Attention model** – This model allows an RNN to pay attention to specific parts of the input that is considered as being important, which improves the performance of the resulting model in practice. By noting $\alpha^{<t,t'>}$ the amount of attention that the output $y^{<t>}$ should pay to the activation $a^{<t'>}$ and $c^{<t>}$ the context at time t , we have:

$$c^{<t>} = \sum_{t'} \alpha^{<t,t'>} a^{<t'>} \quad \text{with} \quad \sum_{t'} \alpha^{<t,t'>} = 1$$

Remark: the attention scores are commonly used in image captioning and machine translation.



A cute teddy bear is reading Persian literature

A cute teddy bear is reading Persian literature

□ **Attention weight** – The amount of attention that the output $y^{<t>}$ should pay to the activation $a^{<t,t'>}$ is given by $\alpha^{<t,t'>}$ computed as follows:

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t''=1}^{T_x} \exp(e^{<t,t''>})}$$

Remark: computation complexity is quadratic with respect to T_x .

* * *

Python For Data Science Cheat Sheet

PySpark - SQL Basics

Learn Python for data science interactively at www.DataCamp.com



PySpark & SparkSQL

SparkSQL is Apache Spark's module for working with structured data.



Initializing SparkSession

A SparkSession can be used create DataFrame, register DataFrame as tables, execute SQL over tables, cache tables, and read parquet files.

```
>>> from pyspark.sql import functions as F  
>>> df = spark.createDataFrame([(1, "John", 20),  
>>> (2, "Pete", 25),  
>>> (3, "Sarah", 30)], ["id", "name", "age"])  
>>> df.show()  
+---+---+---+  
| id | name | age |  
+---+---+---+  
|  1 | John  | 20  |  
|  2 | Pete  | 25  |  
|  3 | Sarah | 30  |  
+---+---+---+
```

```
>>> spark = SparkSession.builder.  
>>> .appName("Python Spark SQL basic example")  
>>> .config("spark.some.config.option", "some-value")  
>>> .getOrCreate()
```

Creating DataFrames

From RDDs

```
>>> from pyspark.sql.types import *  
>>> InferSchema  
>>> sc = spark.sparkContext  
>>> lines = sc.textFile("people.txt")  
>>> lines = lines.map(lambda l: l.split(",")  
>>> people = lines.map(lambda p: Row(name=p[0], age=int(p[1])))  
>>> peopleDf = spark.createDataFrame(people)  
>>> Schema  
>>> people = people.map(lambda p: Row(name=p[0], age=int(p[1].strip())))  
>>> schemaString = "name age"  
>>> fields = [StructField("name", StringType(), True),  
>>> StructField("age", IntegerType(), True)]  
>>> schema = StructType(fields)  
>>> peopleDf = spark.createDataFrame(people, schema).show()  
+---+---+  
| name | age |  
+---+---+  
| John | 20 |  
| Pete | 25 |  
| Sarah | 30 |  
+---+---+
```

From Spark Data Sources

JSON

```
>>> df = spark.read.json("customer.json")  
>>> df.show()
```

```
+---+---+---+  
| address | firstName | lastName |  
+---+---+---+  
| [New York,10021,N... | John | Smith | [1322, 55+-1234,po...  
| [New York,10021,N... | 21 | John | Doe | [322, 88+-1234,po...  
+---+---+---+
```

```
>>> df2 = spark.read.load("people.json", format="json")
```

```
>>> df3 = spark.read.load("users.parquet")
```

```
>>> df4 = spark.read.text("people.txt")
```

Inspect Data

```
>>> df.dtypes  
Return df.column names and data types  
>>> df.show()  
Display the content of df  
Return first n rows  
>>> df.first()  
Return the first n rows  
>>> df.take(2)  
Return the first 2 rows  
>>> df.schema  
Return the schema of df
```

Duplicate Values

Queries

Select

```
>>> from pyspark.sql import functions as F  
>>> df.select("firstName").show()  
>>> df.select("firstName", "lastName")  
>>> df.select("firstName",  
>>> "age",  
>>> "phoneMobile")  
>>> df.select("firstName",  
>>> "age",  
>>> "age")  
>>> df.show()
```

Show all entries in firstName column

Show all entries in firstName, age

and type

Show all entries in firstName and age,

Show all entries where age >= 24

Show firstName and oor depending

on age >30

Show firstName if in the given options

Show FirstName and LastName is

TRUE if LastName is like Smith

Show FirstName and LastName

and LastName starts with Sm

Show last names ending in ch

Show substrings of FirstName

Return substrings of firstName

Show age: values are TRUE if between

22 and 24

Sort

```
>>> df.filter(df["age"] > 24).show()  
>>> df.filter(df["age"] > 24).collect()  
>>> df.sort("age", ascending=False).collect()  
>>> df.orderBy(["age", "city"], ascending=[0,1]).collect()
```

Filter entries of age, only keep those

records of which the values are >24

GroupBy

Duplicate Values

GroupBy

Python For Data Science Cheat Sheet

SciPy - Linear Algebra

Learn More Python for Data Science [Interactively](#) at [www.datacamp.com](#)

SciPy

The SciPy library is one of the core packages for scientific computing that provides mathematical algorithms and convenience functions built on the NumPy extension of Python.



Index Tricks

```
>>> import numpy as np  
>>> a = np.array(([1, 2, 3], [4, 5, 6], [(3, 2, 1), (4, 5, 6)]))  
>>> b = np.array([(1, 5, 2, 3), (4, 5, 6, 5)])  
>>> c = np.array([(1, 5, 2, 3), (4, 5, 6, 5), ((3, 2, 1), (4, 5, 6))])
```

Shape Manipulation

| | |
|--------------------------------|---|
| >>> np.mgrid[0:2, 0:2] | Create a dense meshgrid |
| >>> np.ogrid[[0, 1]*2, :1:10j] | Create an open meshgrid |
| >>> np.c_[b, c] | Stack arrays vertically (row-wise) |
| >>> np.hstack((b, c)) | Stack arrays horizontally (column-wise) |
| >>> np.vstack((a, b)) | Stack arrays vertically (row-wise) |
| >>> np.hsplit(c, 2) | Split the array horizontally at the 2nd index |
| >>> np.vsplit(d, 2) | Split the array vertically at the 2nd index |

Polynomials

```
>>> from numpy import poly1d  
>>> p = poly1d([3, 4, 5])
```

Vectorizing Functions

```
>>> def myfunc(a):  
...     if a < 0:  
...         return a/2  
...     else:  
...         return a/2  
>>> np.vectorize(myfunc)
```

Type Handling

```
>>> np.real(c)  
>>> np.imag(c)  
>>> np.real_if_close(c, tol=1000)  
>>> np.cast['f'](np.pi)
```

Other Useful Functions

```
>>> np.angle(b, deg=True)  
>>> g = np.linspace(0, np.pi, num=5)  
>>> g[3:] += np.pi  
>>> np.unwrap(g)  
>>> np.logspace(0, 10, 3)  
>>> np.select([c<4], [c**2])  
>>> misc.factorial(a)  
>>> misc.comb(10, 3, exact=True)  
>>> misc.central_diff_weights(3)  
>>> misc.derivative(myfunc, 1.0)
```

Linear Algebra

You'll use the linalg and sparse modules. Note that `scipy.linalg` contains and expands on `numpy.linalg`.

Also see [NumPy](#)

Creating Matrices

```
>>> from scipy import linalg, sparse  
>>> A = np.matrix(np.random.random((2,2)))  
>>> B = np.asmatrix(B)  
>>> C = np.mat(np.random.random((10,5)))  
>>> D = np.mat([[3,4], [5,6]])
```

Basic Matrix Routines

| | |
|-------------------------|---------------------------------------|
| Inverse | <code>linalg.inv(A)</code> |
| A.T | <code>linalg.inv(A)</code> |
| Rank | <code>np.linalg.matrix_rank(C)</code> |
| Determinant | <code>linalg.det(A)</code> |
| Solving linear problems | <code>linalg.solve(A, b)</code> |
| Generalized inverse | <code>linalg.pinv(C)</code> |
| Polynomials | <code>linalg.polyv(C)</code> |
| Factorials | <code>linalg.pinv2(C)</code> |

Multiplication

```
>>> np.add(A, D)  
>>> np.subtract(A, D)  
>>> np.divide(A, D)  
>>> np.multiply(D, A)  
>>> np.dot(A, D)  
>>> np.vdot(A, D)  
>>> np.inner(A, D)  
>>> np.outer(B, D)  
>>> np.tensordot(A, D)  
>>> np.kron(A, D)
```

Exponential Functions

```
>>> linalg.expm(A)  
>>> linalg.expm2(A)  
>>> linalg.expm3(D)  
>>> linalg.logm(A)
```

Logarithmic Functions

```
>>> linalg.slogm(B)  
>>> linalg.sinh(D)  
>>> linalg.coshm(D)  
>>> linalg.tanhm(A)  
>>> linalg.tanhm(A)
```

Trigonometric Functions

```
>>> linalg.cosm(B)  
>>> linalg.tanm(B)  
>>> linalg.sinm(D)  
>>> linalg.cosem(D)  
>>> linalg.tanhn(A)  
>>> linalg.tanhn(A)
```

Matrix Sign Function

```
>>> np.signm(A)  
>>> linalg.sqrtm(A)  
>>> linalg.sqrtm(A)
```

Matrix Square Root

```
>>> linalg.sqrtm(A)  
>>> linalg.sqrtm(A)
```

Arbitrary Functions

```
>>> linalg.fsum(A, lambda x: x*x)
```

Decompositions

```
>>> la, v = linalg.eig(A)  
>>> l1, l2 = la  
>>> v[:, 0]  
>>> linalg.eigvals(A)
```

Addition

```
>>> np.add(A, D)
```

Subtraction

```
>>> np.subtract(A, D)
```

Division

```
>>> np.divide(A, D)
```

Multiplication

```
>>> np.multiply(D, A)
```

Dot product

```
>>> np.vdot(A, D)
```

Vector dot product

```
>>> np.inner(A, D)
```

Inner product

```
>>> np.outer(B, D)
```

Outer product

```
>>> np.tensordot(A, D)
```

Tensor product

```
>>> np.kron(A, D)
```

Kronecker product

```
>>> np.expm(A)
```

Matrix exponential

(Taylor Series)

Matrix decomposition

```
>>> np.linalg.eig(B)
```

Matrix logarithm

```
>>> np.linalg.logm(D)
```

Matrix sine

```
>>> np.linalg.sinm(B)
```

Matrix tangent

```
>>> np.linalg.tanm(B)
```

Hyperbolic matrix sine

```
>>> np.linalg.cosem(D)
```

Hyperbolic matrix cosine

```
>>> np.linalg.coshm(D)
```

Hyperbolic matrix tangent

```
>>> np.linalg.tanhm(A)
```

Matrix sign function

```
>>> np.linalg.sqrtm(A)
```

Matrix square root

```
>>> np.linalg.sqrtm(A)
```

Evaluate matrix function

```
>>> linalg.fsum(A, lambda x: x*x)
```

Also see NumPy



Python For Data Science Cheat Sheet

Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.

Learn Python interactively at www.DataCamp.com



1 Prepare The Data

Also see [Lists & NumPy](#)

1D Data

```
>>> import numpy as np
```

```
>>> x = np.linspace(0, 10, 100)
```

```
>>> y = np.cos(x)
```

```
>>> z = np.sin(x)
```

```
>>> data = 2 * np.random.random((10, 10))
```

```
>>> Y, X = np.meshgrid(np.arange(10), np.arange(10))
```

```
>>> U = -1 + X**2 + Y
```

```
>>> V = 1 + Y**2
```

```
>>> from matplotlib.cbook import get_sample_data
```

```
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

```
>>>
```

2D Data or Images

```
>>> data2 = 3 * np.random.random((10, 10))
```

```
>>> ax.plot(x, y, alpha=0.4)
```

```
>>> ax.plot(x, y, c='k')
```

```
>>> ax.colorbar(im, orientation='horizontal')
```

```
>>> im = ax.imshow(img, cmap='seismic')
```

```
>>>
```

Figure

```
>>> fig = plt.figure()
```

```
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
```

```
>>>
```

Axes

```
>>> import matplotlib.pyplot as plt
```

```
>>> fig.add_subplot(111)
```

```
>>> ax1 = fig.add_subplot(221) # row-col-num
```

```
>>> ax2 = fig.add_subplot(222)
```

```
>>> fig3, axes = plt.subplots(nrows=2, ncols=2)
```

```
>>> fig4, axes2 = plt.subplots(ncols=3)
```

```
>>>
```

3 Plotting Routines

```
>>> fig, ax = plt.subplots()
```

```
>>> line = ax.plot(x, y)
```

```
>>> ax.scatter(x,y)
```

```
>>> axes[0,0].bar([1,2,3],[3,4,5])
```

```
>>> axes[0,0].barh([0.5,1.2,5],[0,1,2])
```

```
>>> axes[1,1].axhline(0.45)
```

```
>>> axes[0,1].axvline(0.65)
```

```
>>> ax.fill(x,y,color='blue')
```

```
>>> ax.fill_between(x,y,color='yellow')
```

```
>>>
```

2D Data or Images

```
>>> fig, ax = plt.subplots()
```

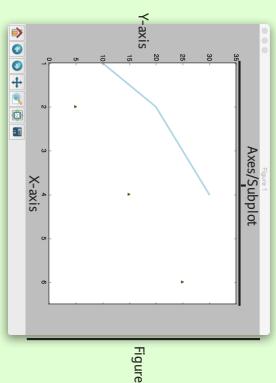
```
>>> im = ax.imshow(img, cmap='gist_earth', interpolation='nearest', vmin=-2, vmax=2)
```

Plot Anatomy & Workflow

Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
- 2 Create plot
- 3 Plot
- 4 Customize plot
- 5 Save plot
- 6 Show plot



4 Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x, x**2, x, x**3)
```

```
>>> ax.plot(x, y, alpha=0.4)
```

```
>>> ax.plot(x, y, c='k')
```

```
>>> ax.colorbar(im, orientation='horizontal')
```

```
>>> im = ax.imshow(img, cmap='seismic')
```

```
>>>
```

Markers

```
>>> fig, ax = plt.subplots()
```

```
>>> ax.scatter(x,y,marker='o')
```

```
>>> ax.plot(x,y,marker='o')
```

```
>>>
```

Linestyles

```
>>> plt.plot(x,y,lw=4.0)
```

```
>>> plt.plot(x,y,lw=1)
```

```
>>> plt.plot(x,y,lw=2,ls='--')
```

```
>>> plt.setp(lines,color='r',linewidth=4.0)
```

```
>>>
```

Text & Annotations

```
>>> ax.text(1, 1,
```

```
'Example Graph',
```

```
style='italic')
```

```
>>> ax.annotate("Sine",
```

```
xy=(8, 0),
```

```
xytext=(10.5, 0),
```

```
arrowprops=dict(arrowstyle="->",
```

```
connectionstyle="arc3"),
```

```
color='red')
```

```
>>>
```

Ticks

```
>>> ax.xaxis.set(ticks=range(1,15),
```

```
label='X-Axis')
```

```
>>> ax.set_xlim(0,10)
```

```
>>> ax.legend(loc='best')
```

```
>>>
```

Limits, Legend & Layouts

```
>>> plt.title(r'$\sigma_1=1.5$', fontsize=20)
```

```
>>> plt.title(r'$\sigma_1=1.5$', font-size=20)
```

```
>>> plt.title(r'$\sigma_1=1.5$')
```

```
>>> plt.title(r'$\sigma_1=1.5$')
```

```
>>>
```

MathText

```
>>> ax.margins(x=0,y=0.1)
```

```
>>> ax.set(xlim=[0,10.5], ylim=[-1.5,1.5])
```

```
>>> fig, ax = plt.subplots()
```

```
>>> ax.set(xlim=[0,10.5], ylim=[-1.5,1.5])
```

```
>>>
```

Legends

```
>>> ax.set(title="An Example Axes",
```

```
ylabel="Y-Axis")
```

```
>>> ax.legend(loc='best')
```

```
>>>
```

Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5,
```

```
hspace=0.3,
```

```
left=0.125,
```

```
right=0.9,
```

```
top=0.9,
```

```
bottom=0.1)
```

```
>>>
```

Axis Spines

```
>>> ax1.spines['top'].set_visible(False)
```

```
>>> ax1.spines['bottom'].set_position(('outward', 10))
```

```
>>> ax1.spines['bottom'].set_color('black')
```

```
>>>
```

Adjusting Spacing

```
>>> fig.tight_layout()
```

```
>>>
```

Fit Subplots

```
>>> fig, ax1 = plt.subplots()
```

```
>>> ax1.set(xlim=[0,10], ylim=[-1.5,1.5])
```

```
>>> ax1.set(xlim=[0,10], ylim=[-1.5,1.5])
```

```
>>>
```

Fit Subplot(s)

```
>>> fig, (ax1, ax2) = plt.subplots(1, 2)
```

```
>>> ax1.set(xlim=[0,10], ylim=[-1.5,1.5])
```

```
>>> ax2.set(xlim=[0,10], ylim=[-1.5,1.5])
```

```
>>>
```

Adjusting Spacing

```
>>> fig, (ax1, ax2) = plt.subplots(1, 2)
```

```
>>> ax1.set(xlim=[0,10], ylim=[-1.5,1.5])
```

```
>>> ax2.set(xlim=[0,10], ylim=[-1.5,1.5])
```

```
>>>
```

Save Plot

```
>>> plt.savefig('foo.png')
```

```
>>>
```

Close & Clear

```
>>> plt.close()
```

```
>>>
```

Save Figures

```
>>> plt.savefig('foo.png')
```

```
>>>
```

Save Transparent Figures

```
>>> plt.savefig('foo.png', transparent=True)
```

```
>>>
```

Close Window

```
>>> plt.close()
```

```
>>>
```

Clear Axis

```
>>> ax.clear()
```

```
>>>
```

Clear Figure

```
>>> plt.clf()
```

```
>>>
```

Close Plot

```
>>> plt.close()
```

```
>>>
```

Label Contour Plot

```
>>> axes2[2].contourf(dat1)
```

```
>>> axes2[2].ax.clabel(cs)
```

```
>>>
```



Python For Data Science Cheat Sheet

3) Renderers & Visual Customizations

Bokeh

Learn Bokeh [interactively](#) at [www.DataCamp.com](#).
taught by Bryan Van de Ven, core contributor.



Plotting With Bokeh

The Python interactive visualization library **Bokeh** enables high-performance visual presentation of large datasets in modern web browsers.



Bokeh's mid-level general purpose `bokeh.plotting` interface is centered around two main components: `data` and `glyphs`.



The basic steps to creating plots with the `bokeh.plotting` interface are:

1. Prepare some data: Python lists, NumPy arrays, Pandas DataFrames and other sequences of values
2. Create a new plot
3. Add renderers for your data, with visual customizations
4. Specify where to generate the output

5. Show or save the results

```
>>> from bokeh.plotting import figure
>>> data = [1, 2, 3, 4]
>>> y = [6, 7, 2, 4, 5]
>>> p = figure(title="simple line example",
    x_axis_label='x',
    y_axis_label='y')
>>> p.line(x, y, legend="Tempo", line_width=2)
>>> output_file("lines.html")  # Step 4
>>> show(p)  # Step 5
```

1) Data

[Also see Lists, NumPy & Pandas](#)

Under the hood, your data is converted to Column Data

Sources. You can also do this manually:

```
>>> import numpy as np
>>> df = pd.DataFrame(np.array([[33, 9, 4, 65, 'US'],
    [32, 4, 4, 66, 'Asia'],
    [21, 4, 4, 105, 'Europe'],
    [1, 4, 4, 105, 'Europe'],
    columns=['mpg', 'cyl', 'hp', 'origin'],
    index=['Toyota', 'Fiat', 'Volkswagen']]))
>>> from bokeh.models import ColumnDataSource(df)
>>> cds_df = ColumnDataSource(df)
```

2) Plotting

```
>>> from bokeh.plotting import figure
>>> layout = row(p1,p2,p3)
>>> p1 = figure(plot_width=300, tools='pan,box_zoom')
>>> p2 = figure(plot_width=300, plot_height=300,
    x_range=(0, 8), y_range=(0, 8))
>>> p3 = figure()
```

Grid Layout

```
>>> from bokeh.layouts import gridplot
>>> row1 = [p1,p2]
>>> row2 = [p3]
>>> layout = gridplot([row1, row2, [p3]])
```

Tabbed Layout

```
>>> from bokeh.models.widgets import Panel, Tabs
>>> tab1 = Panel(child=p1, title="tab1")
>>> tab2 = Panel(child=p2, title="tab2")
>>> tab3 = Panel(child=p3, title="tab3")
>>> layout = Tabs(tabs=[tab1, tab2])
```

Glyphs

```
Scatter Markers
>>> p1.circle(np.array([1, 2, 3]), np.array([3, 2, 1]),
    fill_color='white')
>>> p2.square(np.array([1, 2, 3, 4]), [1, 4, 3],
    color='blue', size=1)

Line Glyphs
>>> p1.line([1, 2, 3, 4], [3, 4, 5, 6], line_width=2)
>>> p2.multi_line([[[1, 2, 3], [5, 6, 7]]],
    pd.DataFrame([[3, 4, 5], [13, 2, 1]]),
    color="blue")
```

Customized Glyphs

```
Selection and Non-Selection Glyphs
>>> p = figure(tools='box_select')
>>> p.circle('mpg', 'cyl', source=cds_df,
    selection_color='red',
    nonselection_alpha=0.1)

Hover Glyphs
>>> p = figure()
>>> p.circle('mpg', 'cyl', source=cds_df,
    tools='hover')
>>> p.add_tools(HoverTool(
    hover_color='blue', tooltips=None, mode='vline')  # Step 1
    , color=dict(cyl='color_mapper'))  # Step 2
    , transform=circle_color_mapper)  # Step 3
```

```
Colormapping
>>> from bokeh.models import CategoricalColorMapper
>>> color_mapper = CategoricalColorMapper(
    factors=['US', 'Asia', 'Europe'],
    palette=['blue', 'red', 'green'])
```

```
>>> p3.circle('mpg', 'cyl', source=cds_df,
    color=dict(cyl='color_mapper'),
    transform=circle_color_mapper,
    legend='Origin')  # Step 4
```

Legend Location

```
Inside Plot Area
>>> p.legend.location = 'bottom_left'
>>> p = figure()  # Step 1
```

Outside Plot Area

```
>>> from bokeh.models import Legend
>>> r1 = p2.asterisk(np.array([[1, 2, 3]]), np.array([3, 2, 1]))
>>> r2 = p2.asterisk(np.array([[4, 5, 6]]), np.array([3, 2, 1]))
>>> legend = Legend(items=[("One", [p1, r1]), ("Two", [p2, r2])],
    location=(0, -30))
>>> p.add_layout(legend, 'right')  # Step 2
```

Legend Orientation

```
>>> p.legend.orientation = "horizontal"
>>> p.legend.location = "vertical"
```

Legend Background & Border

```
>>> p.legend.border_line_color = "navy"
>>> p.legend.background_fill_color = "white"
```

Rows & Columns Layout

```
>>> from bokeh.layouts import row, column
>>> layout = row(p1, p2, p3)
```

5) Show or Save Your Plots

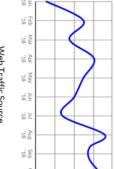
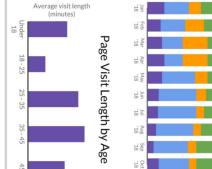
```
>>> show(p1)  # Step 1
>>> save(layout)  # Step 2
```



Exploration and Visualization

The type of dashboard you should use depends on what you'll be using it for.

Common Dashboard Elements

| Type | What is it best for? | Example |
|-------------------|--------------------------------|---|
| Time series | Tracking a value over time |  |
| Stacked bar chart | Tracking composition over time |  |
| Bar chart | Categorical comparison |  |

Popular Dashboard Tools

| Spreadsheets | BI Tools | Customized Tools |
|--------------|----------|------------------|
| Excel | Power BI | R Shiny |
| Sheets | Tableau | d3.js |
| Looker | Looker | |

When You Should Request a Dashboard

- ↻ When you'll use it multiple times
- ⌚ When you'll need the information updated regularly
- == When the request will always be the same

Experimentation and Prediction

Machine Learning

Machine learning is an application of artificial intelligence (AI) that builds algorithms and statistical models to train data to address specific questions without explicit instructions.

| Purpose | Supervised Machine Learning | Unsupervised Machine Learning |
|--|---|-------------------------------|
| Makes predictions from data with labels and features | Makes predictions by clustering data with no labels into categories | |
| Recommendation systems, email subject optimization, churn prediction | Image segmentation, customer segmentation | |

Special Topics in Machine Learning

- **Time Series Forecasting** is a technique for predicting events through a sequence of time and can capture seasonality or periodic events.
- **Natural Language Processing (NLP)** allows computers to process and analyze large amounts of natural language data.
 - Text as input data
 - Word embeddings track the important words in a text
 - Word embeddings create features that group similar words

Deep Learning / Neural Networks enables unsupervised machine learning using data that is unstructured or unlabeled.

Explainable AI is an emerging field in machine learning that applies AI such that results can be easily understood.

Highly accurate predictions

Understandable by humans

Better for "What?"

Better for "Why?"

Python For Data Science Cheat Sheet

Jupyter Notebook

Learn More Python for Data Science Interactively at www.DataCamp.com



Saving/Loading Notebooks

| | | |
|---|----------------------------|--|
| Create new notebook | File → New Notebook | Open an existing notebook |
| Make a copy of the current notebook | File → Open... → ⌘+Shift+N | File → Make a Copy... |
| Save current notebook and record checkpoint | File → Save and Checkpoint | File → Revert to Checkpoint |
| Preview of the printed notebook | Print Preview | File → Download as → - Python notebook - HTML - reStructuredText - LaTeX - PDF |
| Close notebook & stop running any scripts | File → Close and Exit | File → Save and Checkpoint |

Writing Code And Text

Code and text are encapsulated by 3 basic cell types: markdown cells, code cells, and raw NBConvert cells.

Edit Cells

| | | | |
|---|-----------------------------------|--|-----------------------------------|
| Cut currently selected cells to clipboard | Edit → Cut Cells | Copy cells from clipboard to current cursor position | Edit → Copy Cells |
| Paste cells from clipboard above current cell | Edit → Paste Cells | Paste cells from clipboard below current cell | Edit → Paste Cells Below |
| Paste cells from clipboard on top of current cell | Edit → Paste Cells Above | Split up a cell from current cursor position | Cell → Split Cell |
| Revert "Delete Cells" invocation | Cell → Undo Delete Cells | Run current cells down and create a new one above | Cell → Run Cells |
| Merge current cell with the one above | Cell → Merge Cells & Replace | Run all cells above the current cell | Cell → Run All Above |
| Move current cell up | Cell → Move Cell Up | Run all cells below the current cell | Cell → Run All Below |
| Adjust metadata underlying the current notebook | Cell → Edit Notebook Metadata | Change the cell type of current cell | Cell → Cell Type |
| Remove cell attachments | Cell → Remove Cell Attachments | Toggle, toggle scrolling and clear all output | Cell → Current Outputs |
| Paste attachments of current cell | Cell → Paste Cell Attachments | Run current cells down and create a new one below | Cell → All Output |
| Insert Cells | Cell → Insert Cell | Run current cells down and create a new one above | Cell → Run Cells and Insert Below |
| Add new cell above the current one | Cell → Insert Cell Above | Run all cells above the current cell | Cell → Run All |
| | Cell → Kernel → Insert Cell Below | Run all cells below the current cell | Cell → Run All Below |

Executing Cells

| | | | |
|---|------------------------|---|-------------------------------------|
| Run selected cell(s) | Kernel → Run Cells | Run current cells down and create a new one above | Kernel → Run Cells and Insert Below |
| Run current cells down and create a new one above | Kernel → Run Cells | Run all cells above the current cell | Kernel → Run All Above |
| Run all cells below the current cell | Kernel → Run All Below | Run all cells below the current cell | Kernel → Run All Below |
| Change the cell type of current cell | Kernel → Cell Type | Toggle, toggle scrolling and clear all output | Kernel → Current Outputs |
| Toggle, toggle scrolling and clear all output | Kernel → All Output | Run current cells down and create a new one below | Kernel → Run Cells and Insert Below |

Command Mode:

| | | | |
|---------|-------------------------------|-----------------------------------|---|
| In []: | File → Run Cell | File → Kernel → Interrupt | File → Kernel → Clear All Output |
| | File → Kernel → Restart | File → Kernel → Restart & Run All | File → Kernel → Connect back to a remote notebook |
| | File → Kernel → Shutdown | File → Kernel → Reconnect | File → Kernel → Run Other Installed Kernels |
| | File → Kernel → Charge Kernel | | |
| | | | |

Edit Mode:

| | | | |
|---------|------------------------|----------------------|-----------------------|
| In []: | File → Edit Cell | File → Edit Kernel | File → Edit Widgets |
| | File → View Cell | File → View Kernel | File → View Widgets |
| | File → Insert Cell | File → Insert Kernel | File → Insert Widgets |
| | File → Kernel → Kernel | File → Kernel → Help | File → Kernel → Help |
| | | | |

View Cells

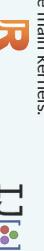
| | | | |
|---|--------------------|---------------------------|-----------------------|
| Toggle display of Jupyter logo and filename | View → Logo | Toggle display of toolbar | View → Toolbar |
| | View → Insert Cell | View → Insert Kernel | View → Insert Widgets |
| | View → Kernel | View → Help | View → Help |
| | | | |
| | | | |

Asking For Help

| | | | |
|------------------------|----------------------------|---|--------------------------------|
| Walk through a UI tour | Help → Help | List of built-in keyboard shortcuts | Help → Keyboard Shortcuts |
| | Help → User Interface Tour | | Help → Edit Keyboard Shortcuts |
| | | Notebook help topics | |
| | | Description of markdown available in notebook | |
| | | Information on unofficial Jupyter notebooks | |
| | | Notebook extensions | |
| | | Information on Python help topics | |
| | | Information on Scipy help topics | |
| | | About Jupyter Notebook | |

Working with Different Programming Languages

Kernels provide computation and communication with front-end interfaces like the notebooks. There are three main kernels:



IPython

| | | | |
|--------------------------------|----------------------------|-----------------------------------|---------------------------|
| Restart kernel | Kernel → Restart | Interrupt kernel | Kernel → Interrupt |
| All cells | Kernel → Interrupt | Clear all output | Kernel → Clear All Output |
| Restart kernel & run all cells | Kernel → Restart & Run All | Connect back to a remote notebook | Kernel → Connect Back |
| Restart kernel & run all cells | Kernel → Restart & Run All | Run other installed kernels | Kernel → Run Other |
| Restart kernel & run all cells | Kernel → Restart & Run All | | |



Julia

| | | | |
|---|---------------------------------|--|--------------------------------------|
| Download serialized state of all widget models in use | Widgets → Download Widget State | Save notebook models in use with interactive widgets | Widgets → Save Notebook with Widgets |
| | | | |
| | | | |
| | | | |
| | | | |



R

| | | | |
|--|--------------------------------------|-----------------------|-------------------------|
| Save notebook models in use with interactive widgets | Widgets → Save Notebook with Widgets | Embed current widgets | Widgets → Embed Widgets |
| | | | |
| | | | |
| | | | |
| | | | |

Widgets

Notebook widgets provide the ability to visualize and control changes in your data, often as a control like a slider, textbox, etc.

Widgets

You can use them to build interactive GUIs for your notebooks or to synchronize stateful and stateless information between Python and JavaScript.



DataCamp

Python for Data Science Cheat Sheet

spacy

Learn more Python for data science interactively at www.datacamp.com



About spaCy

spaCy is a free, open-source library for advanced Natural Language Processing (NLP) in Python. It's designed specifically for production use and helps you build applications that process and "understand" large volumes of text. Documentation: [spacy.io](#)

```
$ pip install spacy
```

```
import spacy
```

Statistical models

Download statistical models
Predict part-of-speech tags, dependency labels, named entities and more. See here for available models: [spacy.io/models](#)

```
$ python -m spacy download en_core_web_sm
```

Check that your installed models are up to date

```
$ python -m spacy validate
```

Loading statistical models

```
import spacy  
# Load the installed model "en_core_web_sm"  
nlp = spacy.load("en_core_web_sm")
```

Documents and tokens

Processing text

Processing text with the `nlp` object returns a `Doc` object that holds all information about the tokens, their linguistic features and their relationships

```
doc = nlp("This is a text")
```

Accessing token attributes

```
doc = nlp("This is a text")  
# Token texts  
[token.text for token in doc]  
#[ 'This', 'is', 'a', 'text' ]
```

Spans

Accessing spans

Span indices are exclusive. So `doc[2:4]` is a span starting at token 2, up to – but not including! – token 4.

```
doc = nlp("I live in New York")  
span = doc[2:4]  
span.text
```

```
# 'a text'
```

Creating a span manually

```
# Import the Span object  
from spacy.tokens import Span  
# Create a Doc object  
doc = nlp("I live in New York")  
# Span for "New York" with label GPE (geopolitical)  
span = Span(doc, 3, 5, label="GPE")  
span.text  
# 'New York'
```

Part-of-speech tags

PREDICTED BY STATISTICAL MODEL

```
doc = nlp("This is a text.")  
[token.pos_ for token in doc]  
#[ 'DET', 'VERB', 'DET', 'NOUN', 'PUNCT' ]  
# Fine-grained part-of-speech tags  
[token.tag_ for token in doc]  
#[ 'DT', 'VBD', 'DT', 'NN', '.' ]
```

Linguistic features

Attributes return label IDs. For string labels, use the `token.pos_`.

Syntax iterators

Sentences

USUALLY NEEDS THE DEPENDENCY PARSER

```
doc = nlp("This is a sentence. This is another one.")  
# doc.sents is a generator that yields spans  
[sent.text for sent in doc.sents]  
#[ 'This is a sentence.', 'This is another one.' ]
```

Base noun phrases

NEEDS THE TAGGER AND PARSER

```
doc = nlp("I have a red car")  
# doc.noun_chunks is a generator that yields spans  
[chunk.text for chunk in doc.noun_chunks]  
#[ 'I', 'a red car' ]
```

Label explanations

```
spacy.explain("RB")  
# 'adverb'  
spacy.explain("GPE")  
# 'Countries, cities, states'
```

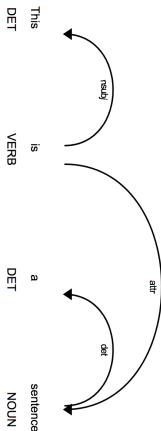
Visualizing

If you're in a Jupyter notebook, use `displacy.render`. Otherwise, use `displacy.sever` to start a web server and show the visualization in your browser.

```
from spacy import displacy
```

Visualize dependencies

```
doc = nlp("This is a sentence")  
displacy.render(doc, style="dep")
```



Visualize named entities

```
doc = nlp("Larry Page founded Google")  
displacy.render(doc, style="ent")
```

Larry Page PERSON founded Google ORG

Word vectors and similarity

To use word vectors, you need to install the larger models ending in `md` or `lg`, for example `en_core_web_lg`.

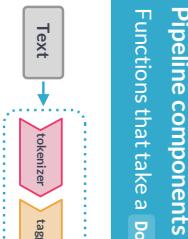
Comparing similarity

```
doc1 = nlp("I like cats")
doc2 = nlp("I like dogs")
# Compare 2 documents
doc1.similarity(doc2)
```

```
# Compare 2 tokens
doc1[1].similarity(doc2[1])
# Compare tokens and spans
doc1[0].similarity(doc2[1:3])
```

Accessing word vectors

```
# Vector as a numpy array
doc = nlp("The sky over New York is blue")
# The L2 norm of the token's vector
doc[2].vector
doc[2].vector_norm
```



Pipeline information

```
nlp = spacy.load("en_core_web_sm")
nlp.pipe_names
# ['tagger', 'parser', 'ner']
nlp.pipeline
# [('tagger', <spacy.pipeline.Tagger>),
# ('parser', <spacy.pipeline.DependencyParser>),
# ('ner', <spacy.pipeline.EntityRecognizer>)]

# Function that modifies the doc and returns it
def custom_component(doc):
    print("Do something to the doc here!")
    return doc

# Add the component first in the pipeline
nlp.add_pipe(custom_component, first=True)
```

Components can be added `first`, `last` (default), or `before` or `after` an existing component.

Extension attributes

Custom attributes that are registered on the global `Doc`, `Token` and `Span` classes and become available as `...`.

Attribute extensions

```
from spacy.tokens import Doc, Token, Span
doc = nlp("The sky over New York is blue")
# The L2 norm of the token's vector
doc[2].vector
doc[2].vector_norm
```

Property extensions

```
# Register custom attribute on Doc class
get_reversed = lambda doc: doc.text[::-1]
Doc.set_extension("reversed", getter=get_reversed)
# Compute value of extension attribute with getter
doc._.reversed
# 'eulb si kxoy weN zevo yks eht'
```

WITH GETTER & SETTER

```
# Register custom attribute on Span class
has_label = lambda span, label: span.label_ == label
Span.set_extension("has_label", method=has_label)
# Compute value of extension attribute with method
doc[3:5].has_label("GPE")
# True
```

Method extensions

CALLABLE METHOD

```
# Register custom attribute on Doc class
get_reversed = lambda doc: doc.text[::-1]
Doc.set_extension("reversed", getter=get_reversed)
# Compute value of extension attribute with getter
doc._.reversed
# True
```

Rule-based matching

Using the matcher

```
# Matcher is initialized with the shared vocab
from spacy.matcher import Matcher
# Each dict represents one token and its attributes
matcher = Matcher(nlp.vocab)
# Add with ID, optional callback and pattern(s)
pattern = [{"LOWER": "new"}, {"LOWER": "york"}]
matcher.add("CITIES", None, pattern)
# Match by calling the matcher on a Doc object
doc = nlp("I live in New York")
matches = matcher(doc)
# Matches are (match_id, start, end) tuples
for match_id, start, end in matches:
    # Get the matched span by slicing the Doc
    span = doc[start:end]
    print(span.text)
```

Rule-based matching

Token patterns

```
# "Love cats", "Loving cats", "Loved cats"
pattern1 = [{"LEMMA": "love", "LOWER": "cats"}]
# "10 people", "twenty people"
pattern2 = [{"LIKE_NUM": True}, {"TEXT": "people"}]
# "book", "a cat", "the sea" (noun + optional article)
pattern3 = [{"POS": "DET", "OP": "?"}, {"POS": "NOUN"}]
```

Operators and quantifiers

Can be added to a token dict as the `"op"` key.

| | |
|---|---|
| ! | Negate pattern and match exactly 0 times. |
| ? | Make pattern optional and match 0 or 1 times. |
| + | Require pattern to match 1 or more times. |

Glossary

Tokenization Segmenting text into words, punctuation etc.

Lemmatization Assigning the base forms of words, for example: "was" → "be" or "rats" → "rat".

Sentence Boundary Finding and segmenting individual sentences.

Part-of-speech (POS) Assigning word types to tokens like verb or noun.

Tagging Assigning syntactic dependency labels, describing the relations between individual tokens, like subject or object.

Named Entity Recognition(NER) Labeling named "real-world" objects, like persons, companies or locations.

Text Classification Assigning categories or labels to a whole document, or parts of a document.

Statistical model Process for making predictions based on examples.

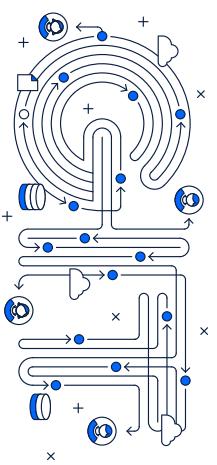
Training Updating a statistical model with new examples.



DataCamp

Learn Python for
data science interactively at
www.datacamp.com

Git Cheat Sheet



| Git Basics | |
|--|---|
| <code>git init <directory></code> | Create empty Git repo in specified directory. Run with no arguments to initialize the current directory as a git repository. |
| <code>git clone <repo></code> | Clone repo located at <repo> onto local machine. Original repo can be located on the local filesystem or on a remote machine via HTTP or SSH. |
| <code>git config user.name <name></code> | Define author name to be used for all commits in current repo. Devs commonly use --global flag to set config options for current user. |
| <code>git add <directory></code> | Stage all changes in <directory> for the next commit. Replace <directory> with a <file> to change a specific file. |
| <code>git commit -m "<message>"</code> | Commit the staged snapshot, but instead of launching a text editor, use <message> as the commit message. |
| <code>git status</code> | List which files are staged, unstaged, and untracked. |
| <code>git log</code> | Display the entire commit history using the default format. For customization see additional options. |
| <code>git diff</code> | Show unstaged changes between your index and working directory. |
| Rewriting Git History | |
| <code>git commit --amend</code> | Replace the last commit with the staged changes and last commit combined. Use with nothing staged to edit the last commit's message. |
| <code>git rebase <base></code> | Rebase the current branch onto <base>. <base> can be a commit ID, a branch name, a tag, or a relative reference to HEAD. |
| <code>git reflog</code> | Show a log of changes to the local repository's HEAD. Add --relative-date flag to show date info or --all to show all refs. |
| Git Branches | |
| <code>git branch</code> | List all of the branches in your repo. Add a <branch> argument to create a new branch with the name <branch>. |
| <code>git checkout -b <branch></code> | Create and check out a new branch named <branch>. Drop the -b flag to checkout an existing branch. |
| <code>git merge <branch></code> | Merge <branch> into the current branch. |
| Remote Repositories | |
| <code>git remote add <name> <url></code> | Create a new connection to a remote repo. After adding a remote, you can use <name> as a shortcut for <url> in other commands. |
| <code>git fetch <remote> <branch></code> | Fetches a specific <branch>, from the repo. Leave off <branch> to fetch all remote refs. |
| <code>git pull <remote></code> | Fetch the specified remote's copy of current branch and immediately merge it into the local copy. |
| <code>git push <remote> <branch></code> | Push the branch to <remote>, along with necessary commits and objects. Creates named branch in the remote repo if it doesn't exist. |
| Undoing Changes | |
| <code>git revert <commit></code> | Create new commit that undoes all of the changes made in <commit>, then apply it to the current branch. |
| <code>git reset <file></code> | Remove <file> from the staging area, but leave the working directory unchanged. This unstages a file without overwriting any changes. |
| <code>git clean -n</code> | Shows which files would be removed from working directory. Use the -f flag in place of the -n flag to execute the clean. |

Additional Options +

| git config | git diff |
|--|--|
| git config --global user.name <name> | git diff HEAD git diff --cached |
| git config --global user.email <email> | Define the author email to be used for all commits by the current user. |
| git config --global alias. <alias-name> | Create shortcut for a Git command. E.g. alias.glog log --graph --oneline --git-command> |
| git config --system core.editor <editor> | Set text editor used by commands for all users on the machine. <editor> arg should be the command that launches the desired editor (e.g., vi). |
| git config --global --edit | Open the global configuration file in a text editor for manual editing. |
| git log | git reset |
| git log --<limit> | Limit number of commits by <limit>. E.g. git log -5 will limit to 5 commits. |
| git log --oneline | Condense each commit to a single line. |
| git log -p | Display the full diff of each commit. |
| git log --stat | Include which files were altered and the relative number of lines that were added or deleted from each of them. |
| git log --author=<pattern> | Search for commits by a particular author. |
| git log --grep=<pattern> | Search for commits with a commit message that matches <pattern>. |
| git rebase | git pull |
| git log --since..<until> | git rebase -i <base> Interactively rebase current branch onto <base>. Launches editor to enter commands for how each commit will be transferred to the new base. |
| git log -- <file> | git pull --rebase <remote> Fetch the remote's copy of current branch and rebases it into the local copy. Uses git rebase instead of merge to integrate the branches. |
| git push | |
| git log --graph --decorate | git push <remote> --force Forces the git push even if it results in a non-fast-forward merge. Do not use the --force flag unless you're absolutely sure you know what you're doing. |
| git log -- <file> | git push <remote> --all Push all of your local branches to the specified remote. |
| git log --graph --decorate | git push <remote> --tags Tags aren't automatically pushed when you push a branch or use the --all flag. The --tags flag sends all of your local tags to the remote repo. |