

MachLe

Summary

Lukas Schöpf

24. Januar 2025

1 Introduction

1.1 Data Types

- Numerical/Quantitive, discret: Countable

Example: Number of rejected Loans, classes taken this semester

- Numerical/Quantitive, Continuos: Interval data

Example: Distance, mg of drug taken, size of a house

- Categorical, ordinal: distinct and can be ordered

Example: Credit score can be {low, medium, high}

- Categorical, nominal: categories cannot be ordered

Example: gender, eye color

1.2 Machine Learning Paradigms

- Unsupervised Learning: Discover and explore structure from unlabelled data

- Supervised Learning: Learn to predict/forecast an output of interest, we know what we want to predict and labelled data is available

1.2.1 Unsupervised Learning

Tasks:

- Dimensionality reduction
- Feature Learning
- Matrix compilation
- Anomaly detection
- Generating data

1.2.2 Supervised Learning

Given a set of features/attributes for some objects and also the output/target value of what we want to predict. The Supervised ML task: Given a new object and its features what would be the output value:

- Regression: Output is a numeric value
- Classification: Output is a categorical value

1.3 Data Preparation/Preprocessing

Data will rarely be in the format and quality needed for analytics and model training and several of these operations will be needed:

- Data integration/consolidation: Collects and merges data from multiple sources into coherent data store
- Data cleaning: removing or modifying incorrect data, identify and reduce noise in data
- Data transformations: normalize, discretize or aggregate the data

- Data reduction: reduce data size by reducing the number of samples or reducing the number of attributes, balance skewed data

Low quality data will result in low quality results

2 Supervised ML, Similarities, kNN & Performance measurements

2.1 Supervised ML

- Goal: Find the best Hypothesis

$$H^* = \arg \min_{H \in \mathcal{H}} \sum_{i=1}^n loss(x_i, y_i)$$

- Loss:

$$loss(x_i, y_i) = (H(x_i) - y_i)^2$$

- Model: $H(x) = a + bx + cx^2 + \dots$

Task of the Learning Algorithm to find best parameters a, b, c (those that minimize the loss)

2.1.1 Overfitting

Model learns training data but doesn't generalize well.

2.1.2 Training and test error

Dataset gets split in Training set and Test set (80%/20%)
 $error_{train}$: Error from trained model on train set
 $error_{test}$: Estimate of the true error (generalization error). Error from trained model on test set.

2.2 kNN

Pros:

- Simple and intuitive
- Multiclass
- Interpretable

Cons:

- Curse of Dimensionality
- Sensitive to noise
- Computationally expensive for large datasets

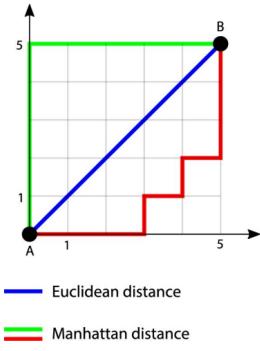
Given training set $X = (x_1, \text{class}_1)(x_2, \text{class}_2)$

Given a new instance $x_?$:

- Find the k-closest examples x_i to $x_?$ in the training set
- Classify $x_?$ based on the majority vote of $\{x_{NN1}, x_{NN2}, \dots\}$

2.2.1 Distance and similarity measurements

Given 2 points $\mathbf{q} = (q_1, q_2, \dots, q_n)$ and $\mathbf{p} = (p_1, p_2, \dots, p_n)$. n is the number of dimensions.



Euclidean Distance

$$d(\mathbf{q}, \mathbf{p}) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

Manhattan Distance

$$d(\mathbf{q}, \mathbf{p}) = \|\mathbf{q} - \mathbf{p}\|_1 = \sum_{i=1}^n |q_i - p_i|$$

2.3 Finding k

$k_{opt} \in \{1, 2, \dots, N\}$ $N = \#$ of training samples

Extreme cases:

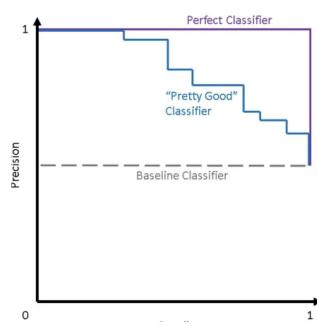
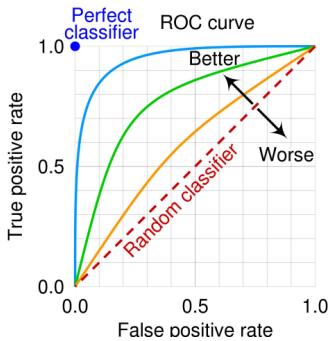
- $k = 1$: kNN = NN
- $k = N$: Majority class

2.4 Performance measures

- Accuracy: $\#corr / \#all = (TP + TN) / (TP + FN + FP + TN)$
- Error: $\#wrong / \#all = 1 - \text{Accuracy}$
- Recall, Sensitivity: $TP / (TP + FN)$, How many relevant samples are correctly detected
- Specificity: $TN / (TN + FP)$
- Precision: $TP / (TP + FP)$, How many detected samples are relevant
- F1 score: $2 \cdot \text{Precision} \cdot \text{Recall} / (\text{Precision} + \text{Recall})$

2.4.1 ROC (Receiver Operating Characteristic) curve

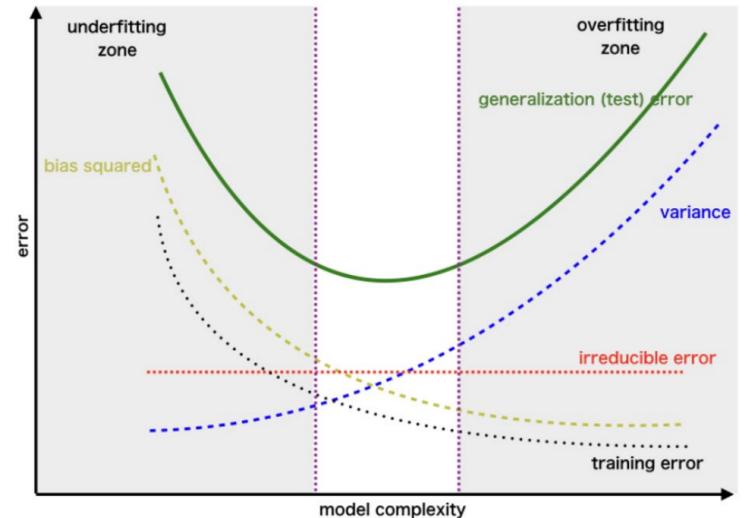
TPR, FPR, to find best threshold.



2.4.2 PRC(Precision-Recall Curve)

Precision, Recall, to find best threshold.

3 Bias-Varinace tradeoff



3.1 No free lunch Theorem

There is no universally best learner (across problems):

3.2 Ockham's Razor

Given 2 models with the same empirical (test) error, the simpler one should be preferred because simplicity is desirable in itself.

3.3 Error sources and bias-variance tradeoff

Different Error sources:

- The model: the best hypothesis is at distance to the true function
- The dataset: different datasets potentially provide different information
- Uncertainty in (X,Y) and its representation:

Partial view of the task: have all relevant features been observed?

Noisy data

Error Decomposit MSE:

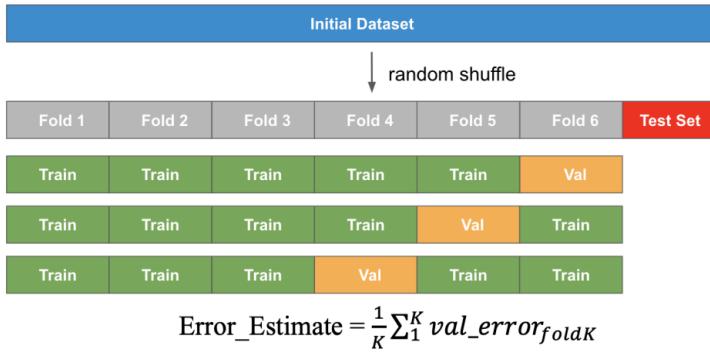
$$E_{MSE} = \text{bias} + \text{variance} + \text{Irreducible error}$$

Bias(systematic error): average predictions deviation from the truth

Variance(dependence on specific sample): Sensitivity of prediction to specific training sample

Irreducible error(random nature of process): due to noise

Generally, for more complex/capable model: bias \downarrow , variance \uparrow . Its a trade-off: Only way to reduce both is to increase the size of the dataset.



3.4 Model selection: Validation score and CV score

k-Fold Cross Validation(CV): We can get a more realistic estimate of the test error using many validation sets (Typically K is 5 to 10)

3.4.1 Error Estimate Summary

In practice:

- Validation score(s) or CV score provide estimates of the test error
- The test error provides an estimate of the true error
- Never use any test data in the model training and model selection process

Which model to choose?

- The one with best validation or CV score
- Use student's t-test to check that an improvement is significant
- Ockham's razor: prefer simpler models in absence of other evidence

Model selection is an empirical science.

3.5 Loss minimization (gradient descent)

Linear Model: $f(x) = w_0 + w_1 x = \hat{y}$

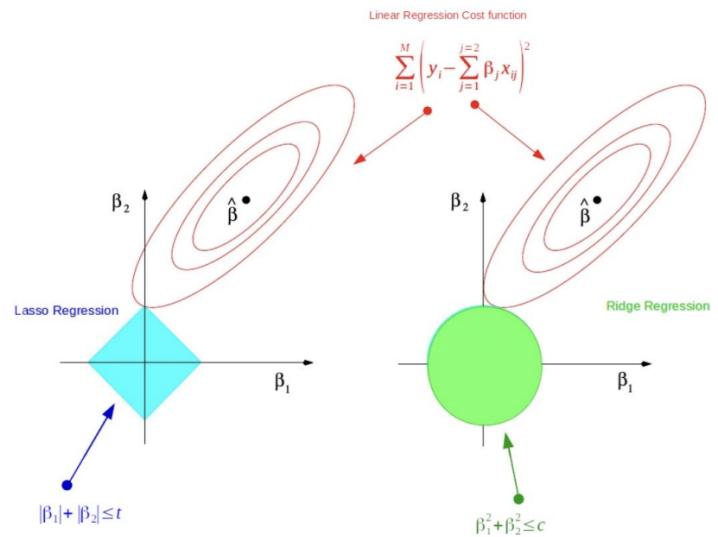
$$RSS = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Use gradient of RSS to find optimal weights for model.

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} - \text{stepsize} \cdot \frac{df(\mathbf{w})}{d\mathbf{w}}$$

The effectiveness of gradient descent depends on the choice of the learning rate (step size):

- Too big, might not reach optimal value
- Too small, it will take a long time to converge



3.6 Regularization

Reduce overfitting of model by penalizing model complexity Tradeoff:

- increase bias
- decrease variance

$$H^* = \arg \min_{H \in \mathcal{H}} \sum_{i=1}^n \mathcal{L}_\theta(x_i, y_i) + \lambda R(H)$$

$R(H)$ is the models complexity. λ controls the model complexity.

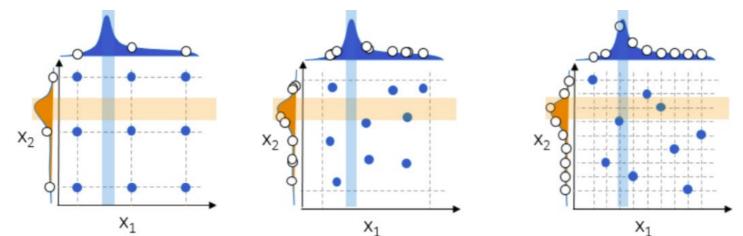
- Lasso:L₁ regularization

$$\arg \min \underbrace{\|y - X\beta\|_2^2}_{\text{Loss}} + \lambda \underbrace{\|\beta\|_1}_{\text{Penalty}}$$

- Ridge:L₂ regularization

$$\arg \min \underbrace{\|y - X\beta\|_2^2}_{\text{Loss}} + \lambda \underbrace{\|\beta\|_2}_{\text{Penalty}}$$

3.7 Hyperparameter tuning



- Grid-search
- Random-search
- {Optimization}

4 Decision Tree, Ensemble Methods & Random Forest

4.1 Decision Tree

- A flow-chart-like tree structure
- Internal node denotes a test on a attribute
- Branch represent an outcome of the test
- Leaf nodes represent class labels or class distribution

Gini Index:

$$I_G = 1 - \sum_{j=1}^c p_j^2$$

Entropy:

$$I_H = - \sum_{j=1}^c p_j \log(p_j)$$

p_j is the proportion of samples that belongs to calss j . Tree construction:

1. Initialization: whole region R_0 (i.e., all given data)
2. Repeat:
 - For each region R_i , for each feature X_j , for each split $R_i = R_{i,l} \cup R_{i,r}$ with respect to feature x_j . Calculate change in impurity score (e.g., gini, entropy, error)
 - Choose best split,i.e., maximum decrease of the impurity score
 - Replace R_i with the two new split regions

Avoiding overfitting:

- Select a proper depth of the tree
- Select a proper minimum number of samples in a leaf to stop further splitting
- Tree pruning - remove split nodes bottom up or top down

All these performed using either a validation set or cross validation!

4.2 Pros/Cons

Pros:

- Easily visualized and interpreted
- No feature normalization or scaling needed
- Works well with mixed feature data types (categorical, continuous)

Cons:

- Easly overfits
- Not robust, high variance

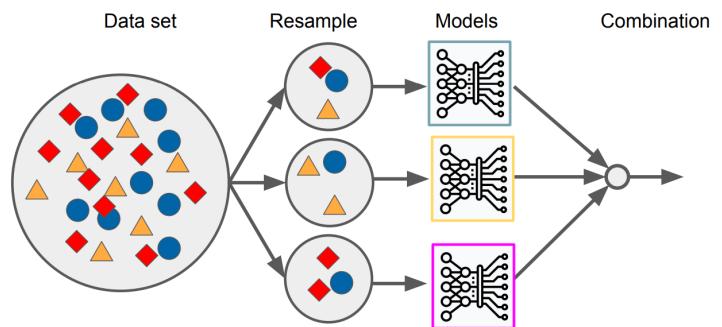
4.3 Ensemble methods

Approach:

- Suppose you have n classifier
- Each classifier has error rate e
- Assume the classifier are independent
- Take Majority vote

The combined result is wrong if $n/2$ classifiers are wrong. According to central limit theorem variance reduces by factor n .

4.3.1 Bagging



To have independent classifier use many independent training sets S_i to train models.

- Variance reduces linearly (sub-linearly in practice because S_i are correlated)
- Bias unchanged (increases slightly in practice)

4.3.2 Boosting

Weight samples. Samples where the model makes mistakes are weighted higher.

Algorithm:

1. Initialization: Train first model on data
2. Repeat:
 - Compute error of the model on each training sample
 - Give higher importance to samples where the model makes mistakes
 - Train next model using importance weighted training samples

In each iteration, introduce a weak model to compensate the shortcoming of the existing strong (= combined) model.

Adaptive Boosting(AdaBoost)

Take a weak learning algorithm and turn it into a strong one by making it focus more on the accurate predictions of difficult cases.

$$F_T(x) = \sum_{t=1}^T f_t(x)$$

In each iteration t :

- A new weak classifier $h(x_i)$ with a coefficient α is added to the existing ones such that the error E_t of the ensemble at iteration t is minimized.

$$E_t = \sum_i [F_{t-1}(x_i) + \alpha_t h(x_i)]$$

- The new weak classifier is trained using a weighted training set where the weight assigned to each sample is identical to the error of the current ensemble classifier on that sample $E(F_{t-1}(x_i))$.

4.3.3 Comparison

No Ensemble:

- complete training set, train one model

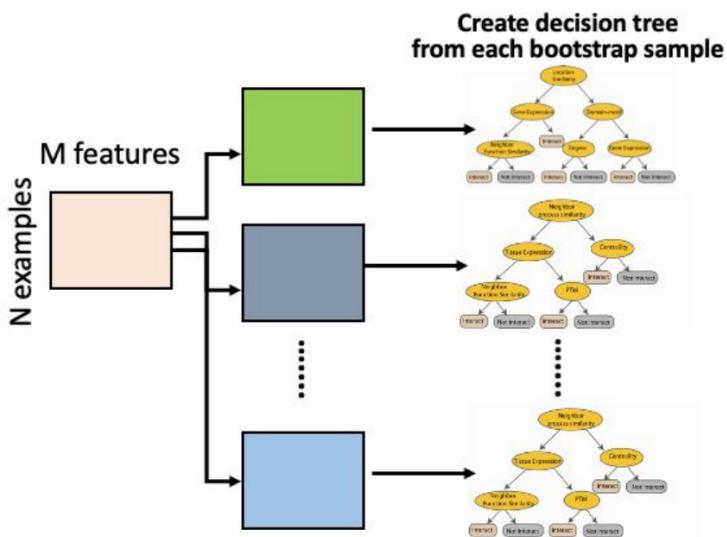
Bagging:

- randomly sample with replacement to obtain different training set
- minimizes variance (usually cannot reduce bias) – fights overfitting
- Computationally efficient (all models can be trained in parallel)

Boosting:

- randomly sample with replacement over weighted data to obtain different training sets
- Minimize bias by adding models to the ensemble – fight underfitting
- Address variance by using simple models with low variance

4.3.4 Random forest

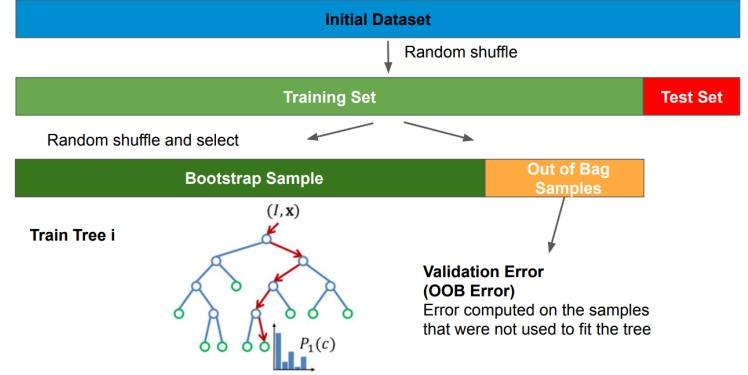


Basic idea:

- Grow many trees in bootstrapped samples of training data
- Minimize bias by growing trees sufficiently deep (overfitting)

- Maximize variance reduction by minimizing correlation between trees by means of bootstrapping data for each tree and sampling variable set at each node
- Reduce variance of noisy but unbiased trees by averaging

Out of Bag Errors (OOB Error)



4.3.5 Summary Random Forrest

Pros:

- Simple - no assumption of the underlying distribution
- OOB error for free
- Many variables, even when they are not relevant for the task at hand or noisy
- Robust against outliers
- Multiclass
- Limit overfitting (trees have to be independent!)
- Unbalanced dataset (subsampling)

5 Probability Recap, Loss Functions, Logistic Regression, Neural Networks Intro

5.1 Probability Basics

Random Variable (RV) x denotes a quantity that is uncertain, discrete or continuous. $p(x = \mathbb{X})$: the probability of variable x being in state \mathbb{X} .

Domain of RV denotes all the values it can take (states it can be in). $\text{dom}(\text{coin}) = \{\text{heads}, \text{tails}\}$

Joint distribution of two RVs x and y takes a particular combination of values and the joint probability density satisfies:

$$\int \int Pr(x, y) \cdot dx dy = 1$$

Marginal distributions $Pr(x)$ and $Pr(y)$ are obtained by:

$$\int Pr(x, y) \cdot dx = Pr(y)$$

$$\int Pr(x, y) \cdot dy = Pr(x)$$

6.2 Training

Repeat:

- Choose a training sample
- Forward pass: Compute the prediction
- Backward pass: If error > 0, update weights

Adjust weights by Gradient descent

$$w_i = w_i - \alpha \frac{\partial J(w_i)}{\partial w_i}$$

Training Terms:

- Epoch = a forward pass and backward pass complete for all the training examples
- Batch size = the number of training samples in a Batch
- Iteration = forward and backward pass each using a Batch
- Iterations per Epoch = #training data /size of Batch

Avoid overfitting with:

- Regularization (Lasso, Ridge)
- Dropout, turn off some neurons during training
- Batch normalization
- Early stopping

6.3 Convolutional Neural Networks (CNN)

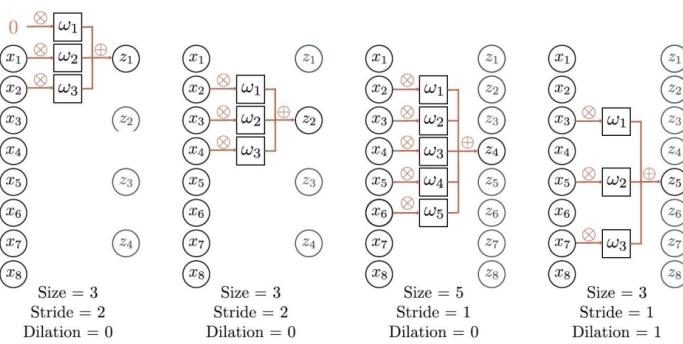
Idea: nearby pixels in an image are correlated - using shared parameters across whole input. A convolution operation is determined:

- stride (kernel shift)
- kernel size (typically odd)
- dilation (number of zero weights in kernel)

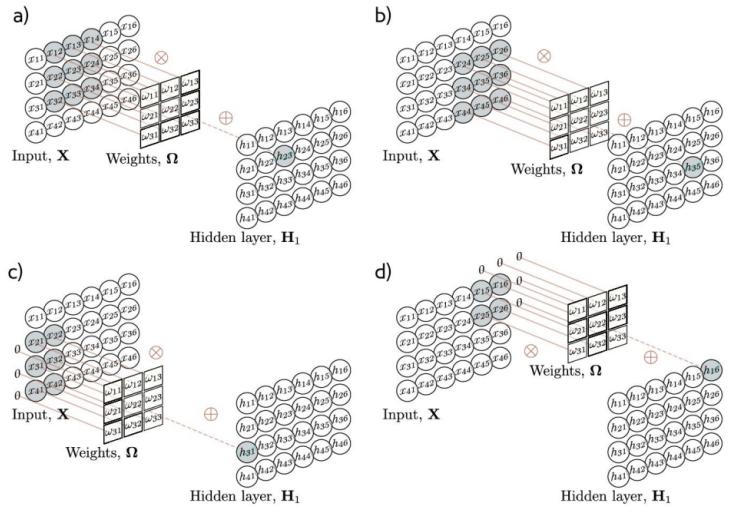
Convolutional Layer: apply convolution, add bias, apply activation function.

Channels(feature maps, activation maps): multiple convolutions applied to the input in parallel.

6.3.1 CNN 1D



1D convolution is a weighted sum of nearby inputs.

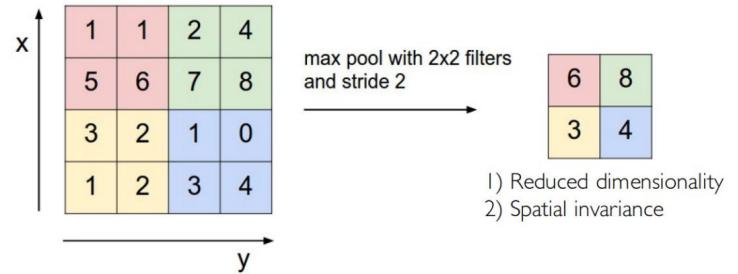


6.3.2 CNN 2D

Using 2D kernels, plus a bias and a activation function.

6.3.3 Pooling

Reduces Dimensions by min,max or average pooling a layer.



6.3.4 Transpose Convolution

Upsampling: Each input contributed multiple times to the output. Useful when output is an image.

7 Feature Engineering

7.1 Data Preparation and Exploration (EDA)

Through exploratory data analysis (EDA) we can often:

- discover important anomalies
- identify limitations in the collection process
- and better inform subsequent goal oriented analysis

Prepare, analyze and visualize data by using:

- Histogramm, QQ-Plot
- Scatter-Matrix
- countplots, contingency tables
- Heatmaps

7.2 Feature Preparation/generation

- **Data Cleaning:** Homogenize missing values and different types of the same feature, fix input errors, types, etc.
- **Aggregation|Pivoting:** Necessary when the entity to model is an aggregation from the provided data.
- **Imputation** of missing values. Strategies: mean, median, mode, using a model
- **Binarization:** Transform discrete or continuous numeric features into binary features
- **Binning** (fixed width, adaptive quantile binning): Split numerical values into bins and encode with a bin ID
- **Transformation**(eg. Log-Transform, BoxCox): Compress the range of large numbers or expand the range of small numbers.
- **Scaling** and **normalizing** (min/max,z-score): Scale numerical variables into a certain range.
- **Generate Interaction features**

7.3 Feature selection

There are general two reasons why feature selection is used:

1. Reducing the number of features, to reduce overfitting and improve the generalization of models
2. To gain a better understanding of the features and their relationship to the response variables

These two goals are often at odds with each other.

7.3.1 Univariate Feature Selection

- Based on univariate statistical tests.Examined for each feature individually
- Good for gaining a better understanding
- For regression: f_regression, mutual_info_regression
- For classification: chi2, f_classif, mutual_info_classif

Pearson Correlation Coefficient:

$$\rho_{X_i, Y} = \frac{\text{cov}[X_i, Y]}{\sigma_{X_i} \cdot \sigma_Y} = \frac{\mathbb{E}[X_i - \bar{X}_i] \cdot \mathbb{E}[Y - \bar{Y}]}{\sigma_{X_i} \cdot \sigma_Y}$$

F-Regression test:

$$\frac{TSS - RSS/p}{RSS/(n - p - 1)} \sim F_{p, n-p-1}$$

Mutual Information coefficient:

$$I(X, Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \cdot \log \left(\frac{p(x, y)}{p(x) \cdot p(y)} \right)$$

7.3.2 Feature selection using linear models & regularization

- Lasso feature selection
- Ridge feature selection
- Linear Models: Subset selection
- Tree-based methods

Ridge feature selection

Forces the coefficient values to be spread out more equally. More useful for feature interpretation than Lasso feature selection.

Linear Models

Forward stepwise selection: Let \mathcal{M}_0 denote the null model with no predictors.

- For $k = 0, \dots, p - 1$:

 1. Consider all $(p - k)$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 2. Choose the best among these $(p - k)$ models and call it \mathcal{M}_{k+1} . Here best is defined as having smallest RSS or highest R^2 .

- Select a single best model among $\{\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p\}$ using AIC or BIC or the adjusted R^2 -score.

Metrics

RSS = Residual Sum of Squares

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

TSS = Total Sum of Squares

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

R^2 -score: (% explained variance)

$$R^2 = 1 - \frac{RSS}{TSS}$$

Adjusted R^2 -score

$$R_{adj}^2 = R^2 = 1 - \frac{\frac{RSS}{n-p-1}}{\frac{TSS}{n-1}}$$

AIC(Akaike) and BIC(Bayesian) Information Criteria, \mathcal{L} : Likelihood

$$AIC = -2 \log(\mathcal{L}) + 2p \quad BIC = -2 \log(\mathcal{L}) + 2 \log(n) \cdot p$$

sparse one-hot
encoding of words

aardvark	1	0	0	...	0	0	0
black	0	0	...	1	...	0	0
cat	0	0	...	1	...	0	0
duvet	0	0	...	1	...	0	0
zombie	0	0	0	...	0	0	1



	animal	fuzziness	dangerous	spooky
aardvark	0.97	0.03	0.15	0.04
black	0.07	0.01	0.20	0.95
cat	0.98	0.98	0.45	0.35
duvet	0.01	0.84	0.12	0.02
zombie	0.74	0.05	0.98	0.93

7.4 Text data/Natural Lannguage Processing (NLP)

Text analysis is a major application field for ML algorithms. To feed text to a algorithms it first has to be transformd to a numerical vector by:

- **Tokenizing:** strings and give an integer ID to each possible token
- **Counting:** the occurrences of token in each document
- **Normalizing:** and weighting with diminishing importance tokens that occur in the majority of samples/documents

7.4.1 Tokenization

1. All common seperators, operators, punctuations and non-printable characters are removed(using Regex).
2. Then, stop-words filtering that aims to filter-out the most frequent words performed
3. Finally, stemming and/or lemmatization is applied to obtain the stem of a word that is morphological root by removing the suffixes that present grammatical or lexical information about the word.

Porter Stemming: Simple replacement rules to create word roots.

Lemmatization: seeks to find actual roots for words.(e.g., tome → book)

7.4.2 Word Embeddings

Idea: Map one-hot encoding to dense vectors

- Each word is represented using a low-dimensional, dense, real-valued vector that encodes the meaning of the word such that the words that are closer in the vector space are expected to be similar in meaning
- How to compute them? Can be pre-trained using a very general task on a very large corpus or learnt with the target task

This is called Word to vector(Word2Vec). There are 2 variants:

- continuous bag of words(CBOW)
- Skip-Gram

7.5 Audio data

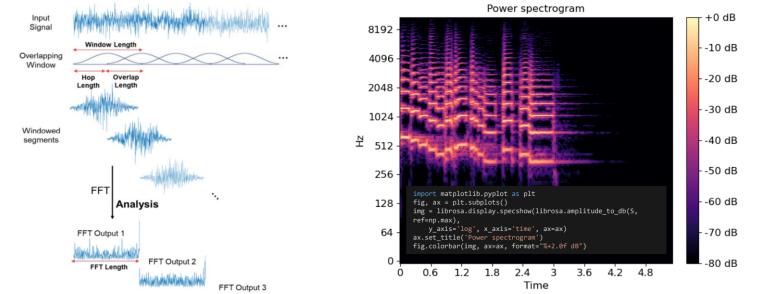
In the time domain the resulting amplitudes, time delays (echos) are captured. The fourier transform enables to decompose a signal into its individual frequencies an their amplitudes and phases.

7.5.1 Linear predictive coding coefficients (LPC)

LPC are parameters derived from audio signals that effectively model the spectral envelope of a sound wave. They capture the articulatory constraints of the vocal tract.

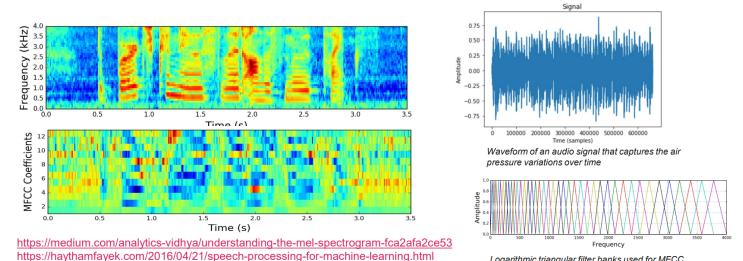
7.5.2 Spectogramm - Short Time Fourier Transform

In many cases we have non periodic signals: frequency content varies over time. Compute FFT on overlapping windowed segments of the signal to obtain the spectogramm.

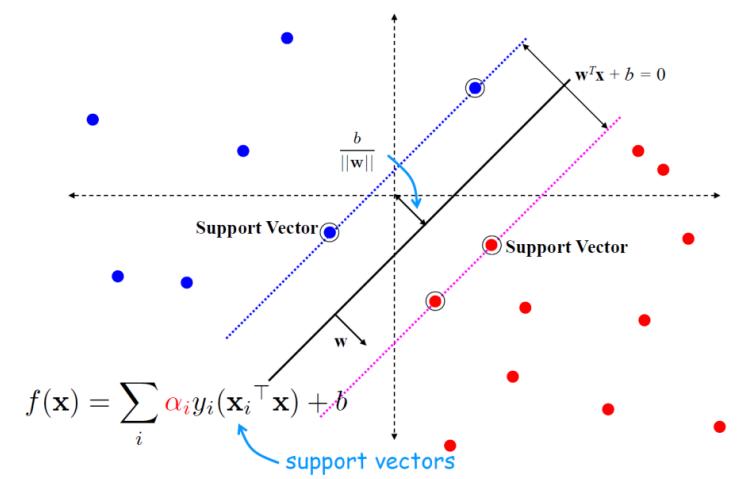


7.5.3 Mel Frequency Cepstral Coefficients (MFCC)

MFCC tries to mimic human hearing. It condenses audio information into fewer coefficients, simplifying data without losing critical information. They are noise robuster than raw spectral features.



8 Support Vector Machines



8.1 Hyperplanes

Hyperplanes can be modified as:

$$H_{w,b} = \{\vec{x} | \vec{w} \cdot \vec{x} + b = 0\}$$

A projection of a point $x \in \mathbb{R}^p$ onto H_w , i.e., the closest point on H_w to x given by:

$$\pi_w(x) = x - \frac{\langle w, x \rangle + b}{\|w\|^2} \cdot w$$

A hyperplane H_w is called separating, if

$$\hat{\gamma}_i = y_i \cdot h(x_i) = y_i \cdot (\langle w, x_i \rangle + b) > 0 \quad (i = 1 \dots n)$$

8.2 maximal margin classifier

The data is called linear separable if there exists a separating hyperplane. In general, if there is one, there are many. We can define our classifier as:

$$f_{w,b}(x) = \text{sign}(\langle w, x \rangle + b)$$

Problem: Not invariant of rescaling. We need some sort of normalization.

We define the geometric margin γ_i of (w, b) with respect to a training sample (x_i, y_i) as:

$$\gamma_i = y_i \cdot \left(\frac{\langle w, x_i \rangle + b}{\|w\|} \right) = \frac{\hat{\gamma}_i}{\|w\|} = y_i d(x_i)$$

The geometric margin with respect of a training set S is smallest of these values in the training set:

$$\gamma = \min_{i=1 \dots n} \gamma_i$$

The geometric margin is invariant to rescaling of the parameters. We want to find a classifier that maximizes the geometric margin for linearly separable dataset.

We will introduce the scaling constraint that the functional margin $\hat{\gamma}$ of (w, b) with respect to the training set S must be 1: $\hat{\gamma} = 1$

$$\begin{aligned} \min_{w,b} & \left\{ \frac{1}{2} \|w\|_2^2 \right\} \\ \text{s.t.} & \quad y_i(\langle w, x_i \rangle + b) \geq 1 \end{aligned}$$

This is a convex optimization problem that could be solved using a Quadratic Programming (QP) code.

8.3 Support Vector Classifier SVC

The hyperplane is only dependent on the vectors closest to it. The closest vectors are called support vectors after their function.

SVM loss function in dual form:

$$-\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \phi(x_i), \phi(x_j) \rangle + \sum_{i=1}^n \alpha_i$$

The Representer Theorem states that the solution w^* of the optimization problem can always be written as a linear combination of the training (or ϕ -transformed training) data:

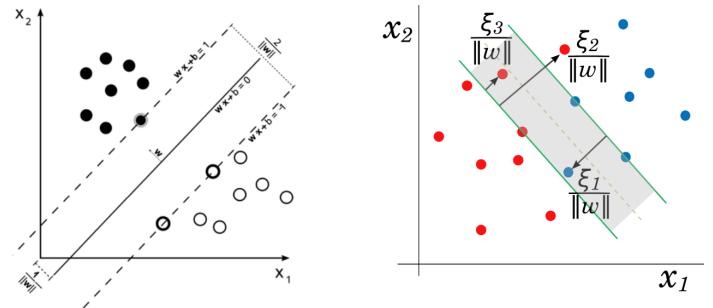
$$w^* = \sum_{i=1}^N \alpha_i y_i x_i \quad w^* = \sum_{i=1}^N \alpha_i y_i \phi(x_i)$$

8.3.1 Handling data that is not linearly separable

1. Apply a feature transform $\phi(x)$
2. Introduce slack variables ξ (tolerate some missclassification)

By increasing the dimensionality of the feature space using mapping ϕ , we can find linearly separating hyperplanes.

8.3.2 Soft margin solution



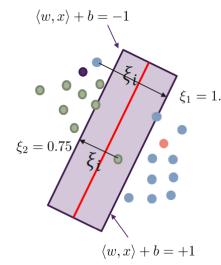
Mathematically, we formulate the trade-off by slack-variables $\xi_i > 0$.

$$\begin{aligned} \min_{w,b} & \left\{ \frac{1}{2} \|w\|_2^2 + C \sum_i \xi_i \right\} \\ \text{s.t.} & \quad y_i(\langle w, \phi(x_i) \rangle + b) \geq 1 - \xi_i \end{aligned}$$

We can fulfill every constraint by choosing ξ_i large enough. The larger ξ_i the larger the objective (that we try minimize). C is a regularization/trade-off parameter:

- small C: allow constraints to be easily ignored. Nearly no penalty for misclassification, large margin
- large C: makes constraints hard to ignore ! smaller margin, strong penalty for misclassification
- C = ∞: recovers a hard margin: all constraints must be fulfilled.

Non separable case example



$$y_1(\langle w, x_1 \rangle + b) \geq 1 - 1.5 = 0.5$$

$$y_2(\langle w, x_2 \rangle + b) \geq 1 - 0.75 = 0.25$$

$$\sum_i \xi_i = 1.5 + 0.75 = 2.25$$

10.4 Mainfold Methods based on similarity

10.4.1 Example of metrics

L_p metric: $d_p(x, y) = \|x - y\|_p = \sqrt[p]{\sum_i |x_i - y_i|^p}$

L_∞ metric: $\|x\|_\infty = \max_i \{|x_i|\}$

L_1 metric: $d_1(x, y) = \|x - y\|_1 = \sum_i |x_i - y_i|$ (manhattan distance)

10.4.2 MDS: Multidimensional scaling

MDS attemps to model similarity or dissimilarity data as distance in geometric spaces. In general, is a technique used for analyzing similarity or dissimilarity data. MDS attemps to find an embedding from the high dimensional objects in I into $y_i \in \mathbb{R}^d$ such that distances d_{ij} are preserved.

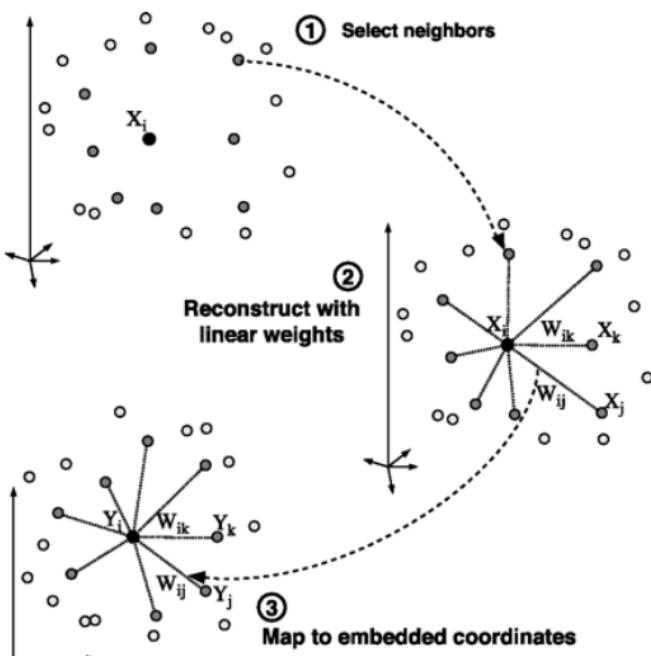
10.4.3 LLE: local linear embedding

LLE describes the local properties of the manifold around a data point x_i by writing the data point as a linear combination (the so-called reconstruction weights w_{ij}) of its k nearest neighbor:

$$x_i \approx \sum_{j=1}^k w_{ij} x_j$$

In the low dimension, LLE attemps to retain the reconstruction weights W as well as possible. Hence, LLE fits a hyperplane through the data point x_i and its nearest neighbors, thereby assuming that the mainfold is locally linear. Finding the low d -dimensional data representation Y amounts to minimizing the cost function \mathcal{L} :

$$\mathcal{L}(w, Y) = \sum_{i=1}^N \|y_i - \sum_{j=1}^k w_{ij} y_j\|^2$$



10.4.4 Isomap: Isometric mapping

10.4.5 t-SNE:t-distributed stochastic neighbor embedding

11 Cluster Analysis

12 Gaussian Mixture Models and EM

13 Reinforcement Learning

14 Generative AI