

### Learning types

**Supervised Learning:** Using annotated (labelled) data. Input Matrix  $X = M(Nr \text{ of training samples}) \times N(Nr \text{ of features})$ . Output Vector  $y = M$ . Goal is to find function  $f(X_{m,:}) \rightarrow y^{(m)}$

**Unsupervised Learning:** No labelled data. Goal is to identify and explain patterns and apply the model to new data => Fuzzy problem state

**Clustering:** Grouping objects in clusters with similar other objects

**Dimensionality Reduction:** Transform data without losing information

**Reinforcement Learning:** Agent observes environment, performs action, gets reward/punishment, adjust policy(strategy) by itself

### Linear Regression

**Hypothesis:**  $h_\theta = \hat{y}^{(m)} = \theta_0 + \theta_1 x^{(m)}$  **Residual:**  $\epsilon^{(m)} = y^{(m)} - \hat{y}^{(m)}$

**Loss:**  $LRSS = \sum_{m=1}^M (y^{(m)} - \hat{y}^{(m)})^2 \approx \sum_{m=1}^M (y^{(m)} - (\theta_0 + \theta_1 x^{(m)}))^2$

**Cost:**  $J(\theta_0, \theta_1) = \frac{1}{2M} \sum_{m=1}^M (y^{(m)} - \hat{y}^{(m)})^2$  **MSE:**  $\frac{1}{M} \sum_{m=1}^M (y^{(m)} - \hat{y}^{(m)})^2$

**Rooted Mean Squared Dev:**  $\sqrt{MSE}$  **MAE:**  $\frac{1}{M} \sum_{m=1}^M |y^{(m)} - \hat{y}^{(m)}|$

**Coefficient of Determination:** Explains the fraction of variance explained by the model.  $R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$ ;  $SS_{res} = \sum (y_i - \hat{y}_i)^2$ ;  $SS_{tot} = \sum (y_i - \mu_y)^2$

**Assumptions:** Linearity (I/O have linear rel.), Independence (residuals/samples are independent), Normality (larger deviations from mean are less likely), Homoscedasticity (Error dist. is the same for all inputs, fan)

sample means

$$\mu_y = \frac{1}{M} \sum_{m=1}^M y^{(m)}$$

$$\mu_x = \frac{1}{M} \sum_{m=1}^M x^{(m)}$$

### Closed Form Equations

$$\theta_0 = \mu_y - \theta_1 \mu_x$$

$$\theta_1 = \frac{\sum_{m=1}^M (y^{(m)} - \mu_y)(x^{(m)} - \mu_x)}{\sum_{m=1}^M (x^{(m)} - \mu_x)^2} = \frac{\tilde{s}_{xy}}{\tilde{s}_x^2}$$

**Multi:**  $\hat{y}^{(m)} = h_\theta(x^{(m)}) = \theta_0 x_0^{(m)} + \theta_1 x_1^{(m)} + \theta_2 x_2^{(m)} + \dots + \theta_N x_N^{(m)} = \theta^T X_{m,:}$

**Dimensions:**  $X: M \times (N + 1)$ ,  $\theta: (N + 1) \times 1$ ,  $y: M \times 1$

### Polynomial Regression

**Hypo:**  $h_\theta(z) = \theta^T z = \sum_{j=0}^n \theta_j z_j$  **Pred:**  $h_\theta(z^{(m)})$

**Regularization:**  $J(\theta) = \frac{1}{2M} [costFn + \lambda \sum_{j=1}^n \theta_j^2]$  Ensure  $\theta$  is not too big

**Hyperparameter:**  $\lambda$ , can be chosen, **tuning**(search) to find "best" one

### Logistic Regression

**Hypo:**  $\hat{y} = h_\theta(x) = g(\theta_0 + \theta_1 x) = 1/(1 + e^{-(\theta_0 + \theta_1 x)})$

**Log-Loss:**  $(h_\theta(x), y) = -y \log(h_\theta(x)) - (1 - y) \log(1 - h_\theta(x))$

Logistic regression (two classes)

$$P(y = 1|x) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$P(y = 1|x) = \frac{1}{1 + e^{-(w^T x)}}$$

$$P(y = 0|x) = 1 - P(y = 1|x) = \frac{e^{-(w^T x)}}{1 + e^{-(w^T x)}}$$

where  $z = w^T x$

### Definitions

**Features:** Individual measurable properties: Independent Variables, Predictors, Attributes or Covariates

**Classification:** Predict a discrete category or class label, fixed nr of output

**Regression:** Predict a continuous numerical value or quantity, find trend

**Instance based:** Pred based on data

**Normal Equation:** For Linear Regress

$$\begin{pmatrix} \epsilon^{(1)} \\ \epsilon^{(2)} \\ \vdots \\ \epsilon^{(M)} \end{pmatrix} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(M)} \end{pmatrix} - \begin{pmatrix} 1 & x^{(1)} \\ 1 & x^{(2)} \\ \vdots & \vdots \\ 1 & x^{(M)} \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix}$$

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

**Residual Plot:** Scatter of pred vs res

**Standard res:** res/st.dev(res)

$$\bar{\epsilon} = \frac{1}{n} \sum_{i=1}^n \epsilon_i \quad s_\epsilon = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (\epsilon_i - \bar{\epsilon})^2}$$

**Model Complexity:** -Bias (Too many assumptions=>underfitting), Variance (+too much attention=>overfitting)

### Performance Metrics

**Accuracy:** corr/all **Precision:** TP/(TP+FP)

**Recall/TPR/Sens:** TP/(TP+FN)

**Confusion Matrix:** TP&TN / PN&PP

**F1:2**  $\frac{prec+recall}{prec+recall} = \frac{2TP}{2TP+FP+FN} \rightarrow 1$  good

**Multi:** Calc for every class and average

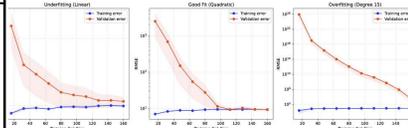
is macro. Micro prec:  $\frac{\sum_{k=1}^K TP_k}{\sum_{k=1}^K (TP_k + FP_k)}$

**K-Fold:** Shuffle. Divide into k-folds. Hold out 1 k for test, train with rest. Repeat with all k. take average of all k errors

**LOOCV:**  $\wedge$  but k=M, good for small data

**Stratified:** Each fold needs same distro

**Learning Curves:** error=train, validation



**Pipeline:** Task\_Def>Data\_Coll+Label>  
**Model** Sel<>Eval> Deploy&Monitor

### Data

**Types:** Categorical=> Nominal (Cannot be ordered, gender) & Ordinal (discrete and ordered), Numerical => Continues & Discrete (can be counted infinite)

**Structured:** Has data model / schema

**Semi-Structured:** contain tags/markers

**Metadata:** Data about data, supportive

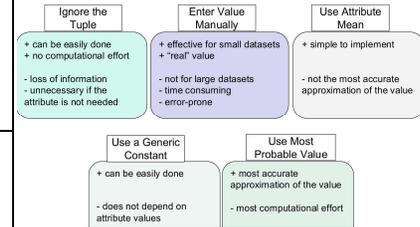
**Cleaning:** Deduplicate, Outlier detection

**Cosine Similarity:** How numerically close are two vectors, only looks at angle not length=>proportional similarity

$$\cos\theta = \frac{A \cdot B}{|A||B|} = \frac{\sum_i^n A_i B_i}{\sqrt{\sum_i^n A_i^2} \sqrt{\sum_i^n B_i^2}}$$

**Levenstein distance:** min nr of edits to transform two words to same

**Missing Values:**



### Probabilities

**Sample Space:** set of all possible outcomes of an experiment. Any subset of the sample space is called an **event** => outcome, or a set of outcomes.

**Probability:** 0-1 of how likely

**Random Variable (RV):** num. quantity from the outcome of random process

**Axioms:** 1. non-negativity 2. Certainty (sum of all=1) 3. Additivity (for mutually exclusive, disjoint events)

**Joint:** A and B / intersection = \*

**Conditional:** that B occurs given that A occurs  $P(B|A) = P(A \cap B) / P(A)$

**Law of total probability:** Calc chance of each unique scenario, add them up weighted by how likely they are

### Gradient Descent

**Repeat:**  $\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ . For 0 nr x

**Batch:**  $\frac{1}{M} \sum_{m=1}^M (h_\theta(x^{(m)}) - y^{(m)}) x^{(m)}$

**SGD:** Shuffle data, for i=1-M, for j=1-N

$$\theta_j = \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

**Learning Rate:** Needs to be in the middle, dynamic. **Decay:**  $\alpha_t = \frac{1}{1 + decay * t} \alpha_0$

**Mini batch:** SGD but batch  $1 < b < M$

### Neural Networks ( $z_j = w_j^T x + w_{0j}$ )

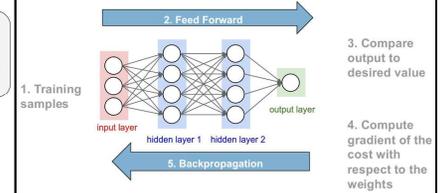
**Softmax Reg:**  $P(y = k|x) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$

$$L_{CE}(\hat{y}^{(m)}, y^{(m)}) = - \sum_{k=1}^K y_k^{(m)} \log \hat{y}_k^{(m)}$$

$$= - \log \hat{y}_c^{(m)} \quad (c \text{ is the correct class})$$

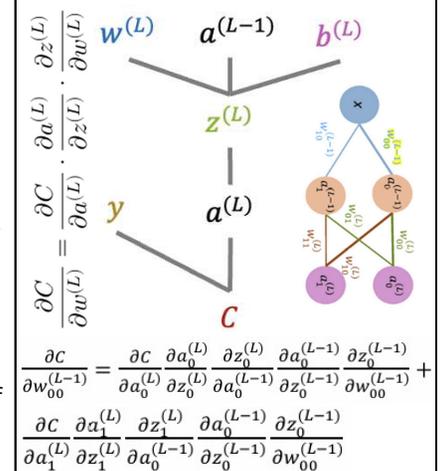
$$= - \log p(y_c^{(m)} = 1|x^{(m)})$$

$$= - \log \frac{\exp(w_c^T x^{(m)})}{\sum_{k=1}^K \exp(w_k^T x^{(m)})}$$



**Dropout:** drop multiple random nodes during each training round

**Backpropagation:** Los: out\*(1-out)



## Support Vector Machines

**Dimensions:** Features **Class Labels:** -1 and +1  
Separate with N-1 dimensional hyperplane  
 $b + w_1x_1 + \dots + w_Nx_N = 0$ , take sign for label

**Margin:** distance of the hyperplane to the closest training samples (gutters)

**Support Vectors:** Datapoints closest to HP

**Distance Point to HP:**  $|w| = \sqrt{w_1^2 + \dots + w_N^2}$

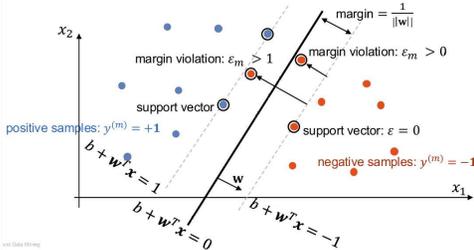
$$r^{(m)} = y^{(m)} \left( \frac{1}{|w|} (w^T x^{(m)} + b) \right)$$

Minimise  $\frac{1}{2} \|w\|^2$  subject to  $y^{(m)}(b + w^T x^{(m)}) \geq 1$  for  $m = 1, \dots, M$

**Normalization:**  $x' = \frac{x - \bar{x}}{s_x}$   $S_x = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N-1}}$

**Soft margin:** Max  $+C \sum_{m=1}^M \epsilon_m$ ; to  $1 - \epsilon_m$

C is hyperparam  $\geq 0$ , controls how many points violate plane (-C=more) => bias-variance?



## Kernel-Trick:

- Linear kernel:**  $K(x^{(i)}, x^{(j)}) = \sum_{n=1}^N x_n^{(i)} x_n^{(j)}$   
- this yields the standard support vector classifier
- Polynomial:**  $K(x^{(i)}, x^{(j)}) = (c + x^{(i)T} x^{(j)})^d = (c + \sum_{n=1}^N x_n^{(i)} x_n^{(j)})^d$   
- Parameters: c and degree d  
- with  $d = 1$  and  $c = 0 \rightarrow$  linear kernel
- Radial basis (RBF):**  
 $K(x^{(i)}, x^{(j)}) = \exp(-\gamma \|x^{(i)} - x^{(j)}\|^2) = \exp\left(-\gamma \sum_{n=1}^N (x_n^{(i)} - x_n^{(j)})^2\right)$   
- Parameter  $\gamma$

## Models Comparison

Type model	+++	---
Generative	Can generate new data samples Better understanding of the data Robust with limited data	More computationally expensive Assumptions of features independence
Discriminative	More accurate for classification Weak assumptions for features independence Better for high-dimensional data	Cannot generate data samples Less interpretable Features-dependent performance

## Decision Trees

**Definitions:** Leaf(end), samples(training instances), value(samples in each class), impure(mix of classes)

**Gini:**  $G(Q_i) = 1 - \sum_{k=1}^K p_{i,k}^2$ ;  $p_{ik} = \frac{1}{M_i} \sum I(y = k)$

How often is random(according to node distro) element incorrectly labelled, max = 0.5

**Weighted Gini:** Weight of childs = goingTo/all.

Then multiply weight\*Gini for each child and add **Information Gain:** Gini(parent)-Weighted Gini^

**Entropy:**  $H(Q_i) = -\sum_{k=1}^K p_{ik} \log_2 p_{ik}$ ; max = 1

**Algorithm:** Separate to maximize information Gain, all possible splits, do again for every split

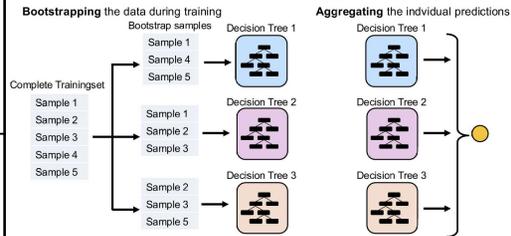
**Regression Trees:** Predict the average numerical value over samples that land in the node, training with CART but for training use MSE

**Regularization:** Pre-Pruning(Min-Sample Pruning: Only split nodes with at least k samples, Max-Depth Pruning: Restrict the maximum depths of the tree), Post-Pruning (Use validation data to decide for each leaf whether it is "reasonable")

**Random Forests:** Train different diverse! trees by bagging and average decision (Ensemble). For classification hard(majority) or soft(class with the highest cumulated probability) voting

**Bagging:** with replacement(pasting without) **Out-of-bag error:** Samples can be left out, rate these elements that were after wrongly classified

**Aggregating the individual predictions**



$$\alpha(x_1) = \frac{3 + 5}{2} = 4$$

$$b(x_1) = \min\left(\frac{6 + 8}{2}, \frac{10 + 12}{2}\right) = 7$$

$$Sil(x_1) = \frac{7 - 4}{7} = \frac{3}{7} = 0.42$$

## Bayes' Theorem

**Bayes Theorem:**  $P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$  Likelihood  $\times$  Prior / Evidence  
disease test that is 99% accurate, 1% of population has disease then. P(Disease)=0.01; P(Pos|Disease)=0.99 so

$$P(\text{Disease}|\text{Positive}) = \frac{P(\text{Positive}|\text{Disease}) \cdot P(\text{Disease})}{P(\text{Positive})}$$

$$P(\text{Positive}) = P(\text{Pos}|\text{Disease}) \cdot P(\text{Disease}) + P(\text{Pos}|\text{No Disease}) \cdot P(\text{No Disease})$$

## Probability Distribution

Normal	$f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	$\mu$ (mean), $\sigma$ (standard deviation) $\mu \in \mathbb{R}, \sigma > 0$
Bernoulli	$P(X = x; p) = p^x (1-p)^{1-x}$ for $x \in \{0, 1\}$	$p$ (probability of success) $0 \leq p \leq 1$
Binomial	$P(X = k; n, p) = \binom{n}{k} p^k (1-p)^{n-k}$ for $k = 0, 1, \dots, n$	$n$ (number of trials), $p$ (probability of success per trial) $n \in \mathbb{N}, 0 \leq p \leq 1$
Beta	$f(x; \alpha, \beta) = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)}$ for $0 \leq x \leq 1$	$\alpha, \beta$ (shape parameters) $\alpha > 0, \beta > 0$
Uniform	$f(x; a, b) = \frac{1}{b-a}$ for $a \leq x \leq b$	$a, b$ (minimum and maximum values) $a < b$

**Maximum Likelihood Estimation (MLE):** Find params to maximize likelihood function. **negative log-likelihood** easier to calc, instead of sum - [log(p)+...] => it becomes a minimization problem

**Maximum A-Posteriori Estimation:**  $\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} \frac{P(X|\theta) \cdot P(\theta)}{P(X)} = P(\theta|X)$  where  $P(X|\theta)$  is the likelihood,  $P(X)$  is the evidence,  $P(\theta)$  is the prior probability

## Clustering (Unsupervised)

**Goals:** Understanding(natural), Identification, reduction(find repr.) and outlier detection **Soft**(Probability for cluster) or **Hard**(unique cluster)

**Similarity:** Jaccard Distance (Intersect/Union), Euclidean  $\|d(p, q)\|^2 = \sum (q_i - p_i)^2$ , Manhattan (straight line,  $d_1(p, q) = \|p - q\|_1 = \sum |p_i - q_i|$ )

## DBSCAN

**Definitions:** minPts (min nr of points (a threshold) clustered together for a region to be considered dense),  $\epsilon$  (distance,  $(x, y) < \epsilon$ ), Core point (has at least minPts points within distance  $\epsilon$  from it), Border point (has at least one core point within distance  $\epsilon$ ), Noise point (is neither) **Properties:** Runtime of  $O(n^2)$  or even  $O(n \log n)$  for indexed data. No need to specify k. arbitrarily shape, noise

## Generative Models

Learn the joint probably for each class  $P(X, Y)$  Of observing class y for features X. Use Bayers to calc  $P(Y|X)$ . **Discriminative** models learn the conditional probability distribution  $(Y|X)$  directly, generative know how data is placed throughout the space. **Inference:** Gen(Highest prob), Disc(decision boundary) **Naïve Bayes:** Assumption all features are conditionally independent => Bernoulli, Multinomial, Gaussian. Works even if not true sometimes

## DBSCAN Noise Algo

Input: Dataset, MinPts,  $\epsilon$   
**1.** Select an unprocessed point P(mark after) **2.** If P is not a Core Point? Classify as noise and go back to **1** **3.** Start new cluster and add P + all neighbours within  $\epsilon$  of it. **4.** For each neighbour that is a core point, recursively add all its neighbours to the cluster, repeat until no more can be added, then go back to **1**

## K-Means

**1.** Init (Random, Forgy(random from training), Partition(Randomly assign a cluster to each point)k centroids **2.** Assign each datapoint to closest k  $f^{(t)}(x_i) = \operatorname{argmin} j \|x_i - c_j^{(t)}\|^2$  **3.** Recalculate centroids as the mean of all points assigned to that cluster=> **2** **4.** Convergence with time complexity  $O(\text{iter} * k * \text{dataPoints} * \text{dim})$

**WCSS:**  $J(C, f) = \sum_{i=1}^n \|x_i - c_{f(x_i)}\|^2$   
**Clusters:** Init has big impact. **Elbow** Run with multiple k values and graph, select k where elbow point is. **Silhouette** Cohesion a is average distance between points in same cluster, Separation b is average distance between point and all points in nearest cluster.  $s = \frac{b-a}{\max(a, b)}$ . High+