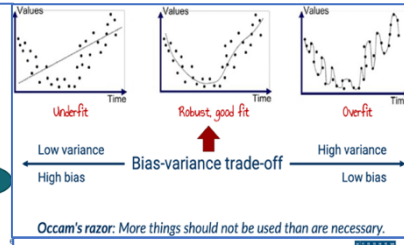
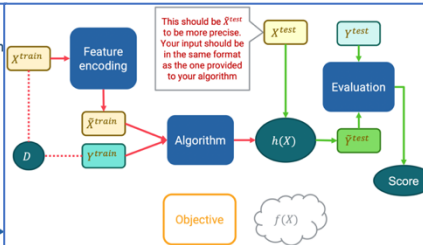
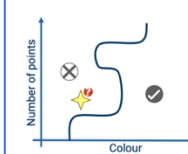


**Lazy Learner**  
Stores the training examples and postpones generalising beyond these data until an explicit request is made at test time.

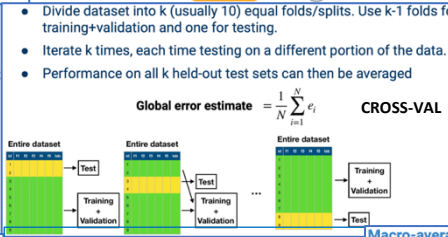


**Eager Learner**  
Constructs a general, explicit description of the target function based on the provided training examples.



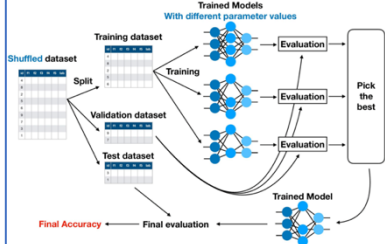
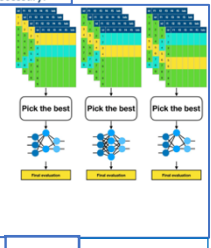
**CASE 1: Plenty of data available -> Held-out test set**  
1) Train algorithm on **training set**  
2) Tune hyperparameters on **validation set**  
3) Estimate generalisation performance using the **test set**

**CASE 2: Limited data available -> Cross-validation**  
1) Separate dataset into **k folds**  
2) Use 1 fold for testing and k-1 folds for **training+validation**  
3) Repeat k times, using each fold as the **test set**  
4) Estimate generalisation performance by **averaging results** across all the test folds



**Option 2: NESTED CROSS-VAL**

- At each cross-validation step separate 1 fold for testing.
- Run an internal cross-validation over the remaining k-1 folds to find the optimal hyperparameters.
- Every fold will still have different hyperparameters, just chosen based on more data.



**Eval metrics**

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Accuracy} = \frac{TP + TN + FP + FN}{N}$$

$$F_1 = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

$$F_1 \text{ score} = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{error rate} = 1 - \text{accuracy}$$

**Macro-averaging:** we are taking the average on the class level. We calculate metrics for each class and average them together.

$$P_{\text{macro}} = \frac{1}{4} \cdot \left( \frac{1}{3} + \frac{1}{1} + \frac{1}{2} + \frac{1}{2} \right) = 0.583$$

**Micro-averaging:** we are taking the average on the item level. For example, adding the TP, FP, TN, FN values for each class before calculating the metrics.

$$P_{\text{micro}} = \frac{1 + 1 + 1 + 1}{(1 + 1 + 1 + 1) + (2 + 0 + 1 + 1)} = 0.5$$

Note: For binary and multi-class classification, micro-averaged P, R and F1 will be just equal to accuracy. This can vary for more complex measures and multi-label classification.

**Macro-avg recall/precision/F1:**  
mean of R/P/F across classes

CM: True/row headers, Imbalanced data, down/upsample maj/min class

The **True error** of the model  $h$  is the probability that it will misclassify a randomly drawn example  $x$  from distribution  $D$ .

$$\text{error}_D(h) \equiv \Pr[f(x) \neq h(x)]$$

The **Sample error** of the model  $h$  based on a data sample  $S$ :

$$\text{error}_S(h) \equiv \frac{1}{N} \sum_{x \in S} \delta(f(x), h(x))$$

$\delta(f(x), h(x)) = 1$  if  $f(x) \neq h(x)$   
 $\delta(f(x), h(x)) = 0$  if  $f(x) = h(x)$

Given a sample  $S$  with  $N \geq 30$ , we can say with  $N\%$  confidence the true error lies in the interval:

Scaling for the desired confidence level  
Both sides of the interval  $\pm Z_{N\%} \sqrt{\text{error}_S(h) \cdot (1 - \text{error}_S(h))}$   
 $n$  = Number of datapoints

The statistical tests tell us if the means of the two sets are significantly different.

There are several statistical tests: Randomisation, T-test, Wilcoxon rank-sum, etc.

Randomisation test: randomly switch some predictions from both models and measure how often the new performance difference is greater or equal than the original difference.

Two-sample T-test: estimate the likelihood that two metrics (e.g. classification error) from different populations are actually different.

Paired T-test: estimating significance over multiple matched results, for example classification error over the same folds in cross-validation

We want to know the true error but we can only measure the sample error.

Distinguish between verifying a hypothesis and exploring the data.

Benjamini & Hochberg (1995) offer an adaptive p-value:

$$p_1 \leq p_2 \leq p_3 \leq \dots \leq p_M$$

$$z_i = 0.05 \cdot \frac{i}{M}$$

3. Significant results are the ones where the p-value is smaller than the critical value.

Entropy is defined as the **average** amount of information:

$$H(X) = - \sum_k P(x_k) \log_2(P(x_k))$$

**Nearest Neighbour classifier**

Classify a test instance to the class label of the nearest training instance (according to some distance metric)

- Non-parametric model
- k Nearest Neighbours

**Algorithm: Decision Tree induction**

- Search for an 'optimal' splitting rule on training data
- Split your dataset according to your chosen splitting rule
- Repeat 1. and 2. on each new splitted subset

**Information Gain**  
 $IG(\text{dataset}, \text{subsets}) = H(\text{dataset}) - \sum_{S \in \text{subsets}} \frac{|S|}{|\text{dataset}|} H(S)$

**Categorical vs. Real-valued attributes**

**Split rule 1:** compute  $H(\text{whole dataset } D)$ ,  $IG(D, \text{feature}_1) \dots IG(D, \text{feature}_n)$ , pick feature w/ highest IG as 1st split. Split rule 2: do same as before, now  $D$  = each subset. For  $IG(D, \text{feature}_1)$ , also need  $H(\text{feature}_1 = \text{cat}_1) \dots H(\text{feature}_1 = \text{cat}_n)$ . To construct tree (find splits) from text: identify the attributes and the categories for each, compute  $IG(D, \text{attribute})$  for all attributes and pick the one w/ highest IG.

**DON'T FORGET NORMALISING WEIGHTS IN IG + compute H(D) based on class labels**

Our loss function (cost function) is the **sum-of-squares**:

$$E = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2$$

**MSE**

$$E = \frac{1}{2} \sum_{i=1}^N (a x^{(i)} + b - y^{(i)})^2 = \frac{1}{2N} \sum_{i=1}^N (Y_i - \hat{Y}_i)^2$$

If  $k > 1$ , pick label w/ highest count among k neighbours. If distKNN, weight each neighbour before summing them

**Manhattan distance ( $L_1$ -norm)**  
 $d(x^{(i)}, x^{(j)}) = \sum_{k=1}^K |x_k^{(i)} - x_k^{(j)}|$

**Euclidean distance ( $L_2$ -norm)**  
 $d(x^{(i)}, x^{(j)}) = \sqrt{\sum_{k=1}^K (x_k^{(i)} - x_k^{(j)})^2}$

**Chebyshev distance ( $L_\infty$ -norm)**  
 $d(x^{(i)}, x^{(j)}) = \max_{k=1}^K |x_k^{(i)} - x_k^{(j)}|$

Others: Mahalanobis distance, Hamming distance, ...

**Definition of derivative:**  $f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$

**Constant rule:**  $\frac{d}{dx}(c) = 0$

**Power rule:**  $\frac{d}{dx}(x^n) = n x^{n-1}$

**Exponential (e):**  $\frac{d}{dx} e^x = e^x$

**Exponential (general):**  $\frac{d}{dx} a^x = a^x \ln(a)$

**Natural log:**  $\frac{d}{dx} (\ln(x)) = \frac{1}{x}$

**Constant multiple rule:**  $\frac{d}{dx}(c f(x)) = c f'(x)$

**Sum rule:**  $\frac{d}{dx}(f(x) + g(x)) = f'(x) + g'(x)$

**Difference rule:**  $\frac{d}{dx}(f(x) - g(x)) = f'(x) - g'(x)$

**Product rule:**  $\frac{d}{dx}(f(x)g(x)) = f(x)g'(x) + g(x)f'(x)$

**Quotient rule:**  $\frac{d}{dx} \left( \frac{f(x)}{g(x)} \right) = \frac{g(x)f'(x) - f(x)g'(x)}{(g(x))^2}$

**Chain rule:**  $\frac{d}{dx} f(g(x)) = f'(g(x))g'(x)$

**Chain rule:**  $\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx}$

**Binary class (binary crossentropy / 1 output sigmoid neuron); multiclass class (cat. crossentropy / >1 softmax neurons); regr (MSE, / one linear neuron); multilabel class (1 binary crossentropy for each sigmoid output neuron / >1 sigmoid neurons)**

**Numerical sol.: gradient descent (iterative)**  
 $\theta_i^{(t+1)} = \theta_i^{(t)} - \alpha \frac{\partial E}{\partial \theta_i^{(t)}}$

**Perceptron: Bin. classif**  
It uses the **threshold function**

$$h(x) = f(W^T x) = \begin{cases} 1 & \text{if } W^T x > 0 \\ 0 & \text{otherwise} \end{cases}$$

Perceptron learning rule:  
 $\theta_i \leftarrow \theta_i + \alpha (y - h(x)) x_i$

Can learn any linearly separable fn

**LINEAR & LOGISTIC REGRESSION**

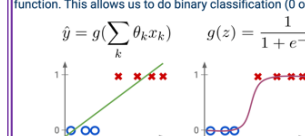
**Simple linear regression: 1 input feature**

**Multiple linear regression: many input features**

e.g. 2 input feat.  $\hat{y}^{(i)} = w_1 x_1^{(i)} + w_2 x_2^{(i)} + b$   
Param updates: use grad of SS loss wrt each param

**Logistic Regression** (Not actually regression)

In logistic regression, we pass the regression output through a logistic function. This allows us to do binary classification (0 or 1).



**General log:**  
 $f(x) = \text{ReLU}(x) = \begin{cases} x & \text{for } x \geq 0 \\ 0 & \text{for } x < 0 \end{cases}$

**1H MLP:**  
 $h \in \mathbb{R}^{H \times 1}$   $W_h \in \mathbb{R}^{K \times H}$   $b_h \in \mathbb{R}^{H \times 1}$   
 $\hat{y} = g_{\hat{y}}(W_{\hat{y}}^T h + b_{\hat{y}})$   $x \in \mathbb{R}^{K \times 1}$  most act. elem-wise except softmax

**Analytical sol.**  
 $\nabla_{\theta} E(\theta) = X^T (X\theta - y) = 0$   
 $\Rightarrow \theta^* = (X^T X)^{-1} X^T y$

Great for directly finding the optimal parameter values. Not so great for large problems: matrix inversion has cubic complexity.

**Maximising the logarithm** is the same as maximising the original product, so let's take the log:

$$L = - \frac{1}{N} \sum_{i=1}^N (y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

Turning this into a loss we can minimise gives us the **binary cross-entropy**:

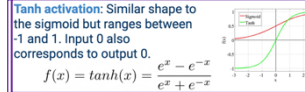
$$L = - \frac{1}{N} \sum_{i=1}^N (y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

This can be generalised to **categorical cross-entropy** for multi-class classification:

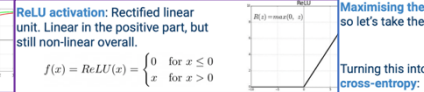
$$L = - \frac{1}{N} \sum_{i=1}^N \sum_{c \in C} y_c^{(i)} \log(\hat{y}_c^{(i)})$$

Where  $C$  is the set of possible classes and  $\hat{y}_c^{(i)}$  is the predicted probability of class  $c$  for datapoint  $i$ .

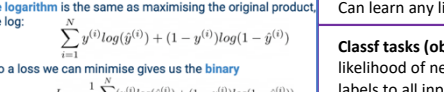
**Tanh activation:** Similar shape to the sigmoid but ranges between -1 and 1. Input 0 also corresponds to output 0.



**ReLU activation:** Rectified linear unit. Linear in the positive part, but still non-linear overall.



**Softmax activation:** Scales the inputs into a probability distribution. The largest input will be large, the rest will be small. All output values sum up to 1.



**CE loss & grad w/ softmax**

$$L = - \frac{1}{N} \sum_{i=1}^N \sum_{c \in C} y_{i,c} \log(\hat{y}_{i,c})$$

$$\frac{\partial L}{\partial z} = \frac{1}{N} (\hat{y} - y)$$

**Batching:** combining the vectors of several datapoints into one matrix to improve speed and reduce noise

**Batching:** combining the vectors of several datapoints into one matrix to improve speed and reduce noise

**Batching:** combining the vectors of several datapoints into one matrix to improve speed and reduce noise

**Batching:** combining the vectors of several datapoints into one matrix to improve speed and reduce noise

**Classf tasks (objective):** maximise the likelihood of network assigning correct labels to all inputs in our dataset

$$\prod_{i=1}^N p(y^{(i)} | x^{(i)}; \theta)$$

**Batching:** combining the vectors of several datapoints into one matrix to improve speed and reduce noise

**Batching:** combining the vectors of several datapoints into one matrix to improve speed and reduce noise

**Don't try to apply this on matrices directly. These are 4D tensors!**

**Scalar-by-vector**

$$\frac{\partial \text{Loss}}{\partial W^{[1]}} = \frac{\partial \text{Loss}}{\partial Z^{[1]}} \cdot \frac{\partial Z^{[1]}}{\partial W^{[1]}}$$

**Vector-by-vector**

$$\frac{\partial \text{Loss}}{\partial W^{[1]}} = \frac{\partial \text{Loss}}{\partial Z^{[1]}} \cdot \frac{\partial Z^{[1]}}{\partial W^{[1]}}$$

**GD (update on N samples, where N=dataset size)**

**Mini-batch GD (update on M samples): most common SGD (update on 1 sample). Adaptive lr: Momentum/Adam/RMSProp/Adagrad/Adadelta**

**Regularisation Techniques**

**K-Means: Summary**

**Step 1: Initialisation**

- Select K generate K random cluster centroids

**Step 2: Assignment**

- Assign each training example to the nearest centroid

**Step 3: Update**

- Update position of each centroid by computing the mean position of all examples assigned to it

**Step 4: Convergence check**

- Stop if position of centroids did not change. Otherwise go to Step 2.

**Better: select K instances at random. Avoids empty clusters.**

**Eucl dist**

**Distance between training example (i) and cluster mean**

**Indicator function 1(a)**

**We have to define k**

- K-means finds a local optimum and is sensitive to the initial centroid positions
- Solutions:
  - Run multiple times, choose model with lowest cost L(O)
  - Use better initialisation: e.g. K-means++

**k-means is not suitable for discovering clusters that are not hyper-ellipsoids (or hyper-spheres).**

**Density Estimation: estimate  $p(x)$  from data**

**Kernel density estimation**

- Compute  $\hat{p}(x)$  by looking at training examples within a kernel function  $H$  (or Parzen window)

**Bandwidth (window size)**

**Kernel function**

**Volume of window**

**Kernel function**

**Increasing the bandwidth  $h$  smooths out  $\hat{p}(x)$**

**Define kernel function  $H$  as a Gaussian: Smoother  $\hat{p}(x)$**

**Gaussian fitting = minimising the negative log likelihood**

- Goal: model captures the probability of generating or observing data  $x$  within an interval
- **Likelihood**: measures the probability of observing data  $x$  from a dataset
- **Negative Log-Likelihood** is often used:
  - Negative: Turn it to a minimisation problem
  - Log: for numerical stability (multiplication of small probabilities  $\rightarrow$  loss)

**Mixture Models**

**Gaussian Mixture Models (most popular)**

**Univariate Gaussian (Normal) distribution**

**Compute  $\hat{p}(x)$  by fitting the parameters  $\hat{\mu}$  and  $\hat{\sigma}^2$  to the training examples**

**Minimize the Bayesian Information Criterion (BIC)**

**Number of parameters**

**Number of examples**

**Number of parameters for 2D Gaussians:  $P_K = 6K - 1$  (2 for mean, 3 for covariance, 1 for mixing proportion)**

**Multivariate Gaussian (Normal) distribution**

**Compute  $\hat{p}(x)$  by fitting the parameters  $\hat{\mu}$  and  $\hat{\Sigma}$  to the training examples**

**Define GA/EA: genotype and phenotype (how go from g to p); fitness fn; selection operator; crossover op; mutation op; hyperparams like pop size**

**GMM-EM**

**Step 1: Initialisation**

- Select K, randomly initialise all parameters

**Step 2: E-Step**

- Compute the responsibilities

**Step 3: M-Step**

- Update the mean
- Update the covariance
- Update the mixing proportions

**Step 4: Convergence check**

- Stop if converged (parameters not changed, or likelihood stagnated). Otherwise go to Step 2.

**K-means**

**Objective function: Minimize average mean squared distance**

**EM-like algorithm**

- Assignment: Assign points to cluster
- Update: Optimise cluster

**GMM-EM**

**Objective function: Maximize log-likelihood**

**EM algorithm**

- E-step: Compute posterior probability of membership
- M-step: Optimise parameters

**Performs hard assignment during Assignment step**

**Performs soft assignment during E-step**

**Assume spherical clusters with equal probability for each cluster**

**Can be used for non-spherical clusters**

**Can generate clusters with different probabilities**

**Genetic algorithm: Optim param list of ints**

- Genotype: binary string of fixed size (ex. 01001010)
- Cross-over: swap portions of the string
- mutation: Random flip of bits

**Genetic programming:**

- Genotype: Program represented as tree (often in LISP)
- Cross-over: Swap portions of trees
- mutation: Change the symbols in the tree

**Evolutionary Strategies: Optim real-val params**

- Genotype: String of floats/doubles
- Cross-over: usually not used
- Mutation: Draw from a Gaussian distribution

**Fitness fn: represent prblm to solve (max. this fn). Gen/Phen: potential sols. Selection op: select parents for next gen (use: Biased Roulette Wheel; Tournament; Elitism).**

**Crossover op: how we combine parents' traits to produce offspring. Mutation op: random var to explore nearby sols.**

**Process: Biased roulette wheel**

- Compute the probability  $p_i$  to select an individual from the population:
- Compute the cumulative probability
- Randomly generate (uniform distribution) a number  $r$  between 0 and 1.
- Select the individual  $x_i$  so that:

**Tournament: Sampling w/ replacement (parents can be selected >1, until we have selected enough parents)**

**Simple version: Tournament**

1. Randomly draw two individuals from the population
2. Select the best out of the two
3. Repeat this until you have enough parents.

**Genetic Algo**

**Initialisation**

**Population**

**Selection**

**Parents**

**Offspring**

**New population**

**Evaluation**

**Selection operator: e.g. tournament**

**Cross-over and Mutation operators**

**Development of the Genotypes into phenotypes**

**Fitness Function**

**Sigma: mut. rate of mutation op**

- Adapting sigma over time (or depending on the situation)
- Concept: Add sigma in the genotype:
  - $x'_j = [x_j, \sigma_j]$
  - $\sigma_j = \sigma_j \exp(\tau_0 \cdot f(0,1))$
  - $y_j = x_j + \sigma_j \cdot f(0,1)$
- $\tau_0$  is the learning rate
- heuristic:  $\tau_0 \propto 1/\sqrt{n}$
- With  $n$  = dimension of the genotype.

**Novelty Search**

**Novelty Archive**

**Novelty**

**Novelty(x) =  $\frac{1}{N} \sum_{i=1}^N d(x, x_i)$**

**Objective: Learning in a single optimisation process a large collection of diverse and high-performing solutions**

**QD-algorithm**

**Stochastic selection**

**Random Mutation**

**Evaluation**

**Tentative addition in the collection**

**Collection**

**Discretises the behavioural descriptor space into a set of cells**

**Addition mechanism:**

- Each new solution goes to the cell corresponding to its BD.
- If the cell is empty the new solution is added to the grid
- If the cell is already occupied, the solution with the best fitness is kept
- Hyper-parameter: size of the cells (or resolution of the grid)
- Advantage: Easy to implement
- Drawback: Density not necessarily uniform

**E.g. QD algo: MAP-Elites = Grid + uni rand select.**

**Coverage = num unique sol / total num sol (or num of cells filled w/ 1 sol in the grid)**

**Diversity = archive size.**

**Performance = max or mean fitness score in archive.**

**Convergence speed = btw diversity and performance metrics. QD-score = sum of fitness scores of all sol in archive (>0).**

**NS w/ local comp: QD algo when used w/ an archive management system**

**Initialisation**

**Population**

**Selection**

**Parents**

**Offspring**

**New population**

**Evaluation**

**Novelty Archive**

**Add most novel individuals in archive**

**Use archive in evaluation**

**Fitness selection**

**Novelty Score**

**Rappel: EuclDist(p,q) =  $\sqrt{(q_x - p_x)^2 + (q_y - p_y)^2}$ , where p and q are R2 vectors. Most novel sol = sol whose BD is furthest away from its neighbours' BD.**