

CSI 4106

Artificial Intelligence

Study Guide



Fall 2024

<h2>Types of Learning</h2> <ul style="list-style-type: none"> Supervised Learning: Each example is accompanied by a label Unsupervised Learning: No feedback is provided to the algorithm Reinforcement Learning: The algorithm receives a reward or a punishment following each action 	<h2>Optimization</h2> <ul style="list-style-type: none"> Until some termination criteria is met: <ul style="list-style-type: none"> Evaluate the loss function, comparing $h(x_i)$ to y_i Make small changes to the weights, in a way that reduces the value of the loss function 	<p>The algorithm is as follows:</p> <p>Repeat until convergence:</p> $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_D)$ <p>for $j \in \{0, \dots, D\}$ (update simultaneously)</p>	<ul style="list-style-type: none"> Some machine learning algorithms are specifically designed to solve binary classification problems <ul style="list-style-type: none"> Logistic regression and SVM are such examples Any multi-class classification problem can be transformed into a binary classification problem One-vs.-All <ul style="list-style-type: none"> A separate binary classifier is trained for each class For each classifier, one class is treated as the positive class, and all other classes are treated as the negative class The final assignment of a class label is made based on the classifier that outputs the highest confidence score for a given input
<h2>Learning Phases</h2> <ul style="list-style-type: none"> Learning (building a model) Inference (using the model) <p>For learning, we first obtain training data, then we convert the data into features vectors. Then we use an algorithm to train a model.</p> <p>For inference, we use new data to make new predictions with the model.</p>	<h2>Gradient Descent (single value)</h2> <ul style="list-style-type: none"> A common loss function for regression problems is the root mean squared error $\sqrt{\frac{1}{N} \sum_{i=1}^N (h(x_i) - y_i)^2}$ <ul style="list-style-type: none"> In practice, minimizing the mean squared error is easier and gives the same result $\frac{1}{N} \sum_{i=1}^N (h(x_i) - y_i)^2$ <ul style="list-style-type: none"> Suppose we have a linear model $h(x_i) = \theta_0 + \theta_1 x_i$ <p>and the loss function</p> $J(\theta_0, \theta_1) = \frac{1}{N} \sum_{i=1}^N (h(x_i) - y_i)^2$ <p>We want to find θ_0 and θ_1 that minimizes J.</p> <ul style="list-style-type: none"> We run gradient descent <ul style="list-style-type: none"> Initialization: θ_0 and θ_1 with random values or zeros Loop: <ul style="list-style-type: none"> Repeat until convergence: $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ for $j \in \{0, 1\}$ α is called the learning rate - size of each step $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ is partial derivative w.r.t. θ_j <p>We would have</p> $J(\theta_0, \theta_1) = \frac{1}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_i - y_i)^2$ $\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{2}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_i - y_i)$ $\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{2}{N} \sum_{i=1}^N x_i (\theta_0 + \theta_1 x_i - y_i)$	<ul style="list-style-type: none"> We assume that the objective function is differentiable <h2>Local vs. Global minimum</h2> <ul style="list-style-type: none"> A function is convex if for every pair of points on the graph of the function, the line connecting these two points lies above or on the graph A convex function has a single minimum <ul style="list-style-type: none"> The loss function for the linear regression (MSE) is convex For functions that are not convex, the gradient descent algorithm converges to a local minimum The loss function generally used with linear or logistic regressions, and SVMs are convex, but not the ones for ANNs 	<h2>Decision Boundary</h2> <ul style="list-style-type: none"> A decision boundary is a boundary that partitions the underlying feature space into regions corresponding to different class labels The data is linearly separable when two classes of data can be perfectly separated by a single linear boundary, such as a line in two-dimensional space or a hyperplane in higher dimensions 2 attributes, linear decision boundary would be a line 2 attributes, non-linear decision boundary would be a non-linear curve 3 attributes, linear decision boundary would be a plane 3 attributes, non-linear decision boundary would be a non-linear surface >3 attributes, linear decision boundary would be a hyperplane >3 attributes, non-linear decision boundary would be a hypersurface Revised definition: a decision boundary is a hypersurface that partitions the underlying feature space into regions corresponding to different class labels
<h2>Supervised Learning</h2> <p>The dataset is a collection of labelled examples</p> <ul style="list-style-type: none"> $\{(x_i, y_i)\}_{i=1}^N$ Each x_i is a feature vector with D dimensions $x_i^{(j)}$ is the value of the feature j of the example i, for $j \in \{1, \dots, D\}$ and $i \in \{1, \dots, N\}$ The label y_i is either a class, taken from a finite list of classes, $\{1, 2, \dots, C\}$, or a real number, or a complex object When the label y_i is a class, taken from a finite list of classes, $\{1, 2, \dots, C\}$, we call the task a classification task When the label y_i is a real number, we call the task a regression task 	<h2>Regression</h2> <ul style="list-style-type: none"> Training data is a collection of labelled examples $\{(x_i, y_i)\}_{i=1}^N$ Each x_i is a feature vector with D dimensions $x_i^{(j)}$ value of feature j of example i, for $j \in \{1, \dots, D\}$ and $i \in \{1, \dots, N\}$ Label y_i is a real number Problem: Given the data set as input, create a model that can be used to predict the value of y for an unseen x 	<h2>Batch Gradient Descent</h2> <ul style="list-style-type: none"> This algorithm is known as batch gradient descent since for each iteration, it processes the "whole batch" of training examples This algorithm might take more time to converge if the features are on different scales The batch gradient descent becomes very slow as the number of training examples increases This is because all the training data is seen at each iteration 	<h2>Logistic Regression</h2> <ul style="list-style-type: none"> Despite its name, logistic regression serves as a classification algorithm rather than a regression technique The labels in logistic regression are binary values, denoted as $y_i \in \{0, 1\}$, making it a binary classification task The primary objective of logistic regression is to determine the probability that a given instance x_i belongs to the positive class, i.e. $y_i = 1$ The standard logistic function is defined as $\sigma: \mathbb{R} \rightarrow (0, 1)$ $\sigma(t) = \frac{1}{1 + e^{-t}}$
<h2>Linear Regression</h2> <ul style="list-style-type: none"> A linear model assumes that the value of the label, \hat{y}_i, can be expressed as a linear combination of the feature values $x_i^{(j)}$: $\hat{y}_i = \theta_0 + \theta_1 x_i^{(1)} + \theta_2 x_i^{(2)} + \dots + \theta_D x_i^{(D)}$ <ul style="list-style-type: none"> Here, θ_j is the j^{th} parameter of the (linear) model, with θ_0 being the bias term/parameter, and $\theta_1, \dots, \theta_D$ being the feature weights Problem: find values for all the model parameters so that the model "best fits" the training data Root mean square error is a common performance measure for regression problems 	<h2>Multivariate Gradient Descent</h2> <ul style="list-style-type: none"> A multivariate linear regression model is defined as $h(x_i) = \theta_0 + \theta_1 x_i^{(1)} + \dots + \theta_D x_i^{(D)}$ $x_i^{(j)}$ = value of feature j in i^{th} example D = number of features The new loss function is $J(\theta_0, \theta_1, \dots, \theta_D) = \frac{1}{N} \sum_{i=1}^N (h(x_i) - y_i)^2$ <p>Its partial derivative</p> $\frac{\partial}{\partial \theta_j} J(\theta) = \frac{2}{N} \sum_{i=1}^N x_i^{(j)} (h(x_i) - y_i)$ <p>where θ, x_i and y_i are vectors, and θx_i is a vector operation</p> <ul style="list-style-type: none"> The vector containing the partial derivatives of J is called the gradient vector $\nabla_{\theta} J(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} J(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_D} J(\theta) \end{pmatrix}$ <ul style="list-style-type: none"> This vector gives the direction of the steepest ascent It gives its name to the gradient descent algorithm $\theta = \theta - \alpha \nabla_{\theta} J(\theta)$	<h2>Stochastic Gradient Descent</h2> <ul style="list-style-type: none"> The stochastic gradient descent algorithm randomly selects one training instance to calculate its gradient This allows it to work with large training sets Its trajectory is not as regular as the batch algorithm <ul style="list-style-type: none"> Because of its bumpy trajectory, it is often better at finding the global minimum when compared to batch Its bumpy trajectory makes it bounce around the local minimum 	<h2>Classification</h2> <ul style="list-style-type: none"> Binary classification is a supervised learning task where the objective is to categorize instances (examples) into one of two discrete classes A multi-class classification task is a task of supervised learning problem where the objective is to categorize instances into one of three or more discrete classes
<h2>Building Blocks of Learning Algorithms</h2> <ul style="list-style-type: none"> A typical learning algorithm comprises the following components: <ul style="list-style-type: none"> A model, often consisting of a set of weights whose values will be "learnt" An objective function Optimization algorithm 	$\frac{1}{N} \sum_{i=1}^N (h(x_i) - y_i)^2$	<h2>Mini-batch Gradient Descent</h2> <ul style="list-style-type: none"> At each step, rather than selecting one training example as SGD does, mini-batch gradient descent randomly selects a small number of training examples to compute the gradients Its trajectory is more regular compared to SGD <ul style="list-style-type: none"> As the size of the mini-batches increases, the algorithm becomes increasingly similar to batch gradient descent, which uses all the examples at each step It can take advantage of the hardware acceleration of matrix operation, particularly with GPUs 	<ul style="list-style-type: none"> An S-shaped curve, such as the standard logistic function, is termed a squashing function because it maps a wide input domain to a constrained output range. Analogous to linear regression, logistic regression computes a weighted sum of the weighted features, expressed as: $\sigma(\theta_0 + \theta_1 x_i^{(1)} + \theta_2 x_i^{(2)} + \dots + \theta_D x_i^{(D)})$ However, using the standard logistic function limits its output to the range $(0, 1)$: $\sigma(\theta_0 + \theta_1 x_i^{(1)} + \theta_2 x_i^{(2)} + \dots + \theta_D x_i^{(D)})$ The logistic regression model, in its vectorized form, is defined as: $h(x_i) = \sigma(\theta x_i) = \frac{1}{1 + e^{-\theta x_i}}$
			<ul style="list-style-type: none"> In logistic regression, the probability of correctly classifying an example increases as its distance from the decision boundary increases The principle holds for both positive and negative classes An example lying on the decision boundary has a 50% probability of belonging to either class

<p>Predictions are made as follows:</p> <ul style="list-style-type: none"> $y_i = 0$, if $h_\theta(x_i) < 0.5$ $y_i = 1$, if $h_\theta(x_i) \geq 0.5$ <p>The values of θ are learned using gradient descent</p> <h3>Feature Engineering</h3> <ul style="list-style-type: none"> The process of creating, transforming, and selecting variables from raw data to improve the performance of machine learning models 	<h3>Recall</h3> <ul style="list-style-type: none"> Also known as sensitivity or true positive rate (TPR) $\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$ <ul style="list-style-type: none"> Proportion of true positive instances correctly identified among all actual positive instances 	<h3>Errors</h3> <ul style="list-style-type: none"> Training Error: Generally tends to be low. Achieved by optimizing learning algorithms to minimize errors through parameter adjustments. Generalization Error: Error rate observed when the model is evaluated on new, unseen data. 									
<h3>Underfitting and Overfitting</h3> <ul style="list-style-type: none"> Underfitting: <ul style="list-style-type: none"> Model is too simple Uninformative features Poor performance on both train and test data Overfitting: <ul style="list-style-type: none"> Model is too complex Too many features Excellent performance on training set, but poor performance on test set 	<h3>F₁ Score</h3> <ul style="list-style-type: none"> Harmonic mean of precision and recall $F_1 \text{ score} = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}} = 2 \times \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{recall}} = \frac{\text{TP}}{\frac{\text{TP} + \text{FP} + \text{FN}}{2}}$	<ul style="list-style-type: none"> Underfitting: <ul style="list-style-type: none"> High training error Model is too simple to capture underlying patterns Poor performance on both training and new data Overfitting: <ul style="list-style-type: none"> Low training error, but high generalization error Model captures noise or irrelevant patterns Poor performance on new, unseen data 									
<h3>Learning Curves</h3> <ul style="list-style-type: none"> We can visualize learning curves to assess our models A learning curve shows the performance of the model on both training and test set Multiple measurements are obtained by repeatedly training the model on larger and larger subsets of the data 	<h3>Micro Performance Metrics</h3> <ul style="list-style-type: none"> Micro performance metrics aggregate the contributions of all classes to compute the average performance metric, such as precision, recall, or F₁ score. This approach treats each individual prediction equally, providing a balanced evaluation by emphasizing the performance on frequent classes 	<h3>Cross-Validation</h3> <ul style="list-style-type: none"> A method used to evaluate and improve the performance of machine learning models It involves partitioning the dataset into multiple subsets, training the model on some subsets while validating it on the remaining ones 									
<h3>Bias / Variance Tradeoff</h3> <ul style="list-style-type: none"> Bias is the error from overly simplistic models. High bias can lead to underfitting Variance is the error from overly complex models. High variance can lead to overfitting Tradeoff: aim for a model that generalizes well to new data 	<h3>Macro Performance Metrics</h3> <ul style="list-style-type: none"> Compute the performance metric independently for each class and then average these metrics. Treats each class equally, regardless of frequency 	<h3>k-Fold Cross-Validation</h3> <ol style="list-style-type: none"> Divide the dataset into k equally sized parts (folds) Training and validation <ul style="list-style-type: none"> For each iteration, one fold is used as the validation set, the remaining k-1 folds are used as the training set Evaluation: The model's performance is evaluated in each iteration, resulting in k performance measures Aggregation: Statistics are calculated based on k performance measures 									
<h3>Confusion Matrix</h3> <table border="1"> <thead> <tr> <th></th> <th>Positive (predicted)</th> <th>Negative (predicted)</th> </tr> </thead> <tbody> <tr> <td>Positive (Actual)</td> <td>True Positive (TP)</td> <td>False Negative (FN)</td> </tr> <tr> <td>Negative (Actual)</td> <td>False Positive (FP)</td> <td>True Negative (TN)</td> </tr> </tbody> </table>		Positive (predicted)	Negative (predicted)	Positive (Actual)	True Positive (TP)	False Negative (FN)	Negative (Actual)	False Positive (FP)	True Negative (TN)	<h3>Precision-Recall Tradeoff</h3> <ul style="list-style-type: none"> As the decision threshold decreases, a higher number of examples are predicted positive, potentially leading the classifier to eventually label all instances as positive Conversely, as decision threshold increases, fewer examples are classified as positive, which may result in the classifier predicting no positive instances at all 	<h3>Properties of Cross-Validation</h3> <ul style="list-style-type: none"> More reliable: <ul style="list-style-type: none"> More reliable estimate of model performance compared to a single train-test split Reduces the variability associated with a single split, leading to a more stable and unbiased evaluation For large values of k, consider the average, variance, and confidence interval Better Generalized <ul style="list-style-type: none"> Helps in assessing how the model generalizes to an independent dataset It ensures that the model's performance is not overly optimistic or pessimistic by averaging results over multiple folds Efficient Use of Data <ul style="list-style-type: none"> Particularly beneficial for small datasets, cross-validation ensures that every data point is used for both training and validation This maximizes the use of available data, leading to more accurate and reliable model training Hyperparameter Tuning
	Positive (predicted)	Negative (predicted)									
Positive (Actual)	True Positive (TP)	False Negative (FN)									
Negative (Actual)	False Positive (FP)	True Negative (TN)									
<h3>Accuracy</h3> <ul style="list-style-type: none"> The accuracy of a model measures how accurate the result is Ratio of correct number of predictions to total number of predictions $\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{\text{TP} + \text{TN}}{\text{N}}$ <p>Accuracy can be misleading in the case of class imbalance, as it disproportionately reflects the performance on the majority class, masking poor performance on the minority class</p>	<h3>ROC Curve</h3> <p>Receiver Operating Characteristics curve</p> <ul style="list-style-type: none"> Curve shows the true positive rate against false positive rate An ideal classifier has TPR close to 1.0 and FPR close to 0.0 $\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$ TPR approaches one when the number of false negative prediction is low $\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$ FPR approaches zero when the number of false positive is low 	<h3>Machine Learning Engineering</h3> <ol style="list-style-type: none"> Gather adequate data Extract features from the raw data <ul style="list-style-type: none"> This process is labor-intensive It necessitates creativity Domain knowledge is highly beneficial 									
<h3>Precision</h3> <ul style="list-style-type: none"> Also known as positive predictive value Proportion of true positive predictions among all positive predictions $\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ <p>An algorithm that makes a small number of high-confidence predictions might achieve a high precision score, but may not be helpful</p>	<h3>Training and Test Set Evaluation</h3> <ul style="list-style-type: none"> Guideline: Typically, allocate 80% of dataset for training and reserve 20% for testing Training Set: Utilized to train model Test Set: Independent subset used exclusively at the final stage to assess the model's performance 	<h3>One-Hot Encoding</h3> <ul style="list-style-type: none"> Some learning algorithms require data to be in numerical form We want to find the best encoding method One-Hot Encoding: Should be preferred for categorical data <ul style="list-style-type: none"> Increases Dimensionality: one-hot encoding increases the dimensionality of feature vectors Avoids Bias: Other encoding methods can introduce biases One-Hot Encoding is a technique that converts categorical variables into a binary vector representation where each category is represented with a single 1 and all other elements as 0 									

<h3>Discretization</h3> <ul style="list-style-type: none"> Grouping ordinal values into discrete categories Binning, bucketing, quantization Advantage: <ul style="list-style-type: none"> Enables the algorithm to learn effectively with fewer training examples Disadvantages: <ul style="list-style-type: none"> Requires domain expertise to define meaningful categories May lack generalizability 	<h3>Using ML for Imputation</h3> <ul style="list-style-type: none"> Predict unknown labels for given examples Can be framed as a supervised learning problem $X_i^{(1)}, X_i^{(2)}, \dots, X_i^{(C-1)}, X_i^{(C)}, \dots, X_i^{(n)}$ Let $\hat{y}_i = X_i^{(C)}$ Training set: use examples where $X_i^{(C)}$ is not missing Method: Train a classifier on this set to predict the missing values Advantages: <ul style="list-style-type: none"> Effectively handle complex relationships and correlations between features Disadvantages: <ul style="list-style-type: none"> Cost-intensive in terms of labor, CPU time, and memory resources 	<h3>Universal Approximation Theorem</h3> <p>A feedforward neural network with a single hidden layer containing a finite number of neurons can approximate any continuous function on a compact subset of \mathbb{R}^n, given appropriate weights and activation functions</p>
<h3>Normalization Scaling</h3> <ul style="list-style-type: none"> Learning algorithms perform optimally when feature values have similar ranges For every sample, we do $\frac{x_i^{(i)} - \min^{(i)}}{\max^{(i)} - \min^{(i)}}$	<h3>Advantages:</h3> <ul style="list-style-type: none"> Transforms each feature to have mean 0 and standard deviation 1 For every sample, we do $\frac{(x_i^{(i)} - \mu^{(i)})}{\sigma^{(i)}}$	<h3>Deep Learning</h3> <ul style="list-style-type: none"> A machine learning technique that can be applied to supervised learning, unsupervised learning, and reinforcement learning Inspired from the structure and function of biological neural networks found in animals Comprises interconnected neurons arranged into layers
<h3>Standardization Scaling</h3> <ul style="list-style-type: none"> Range of values are not bounded 	<h3>Class Imbalance</h3> <ul style="list-style-type: none"> Transforms each feature to have mean 0 and standard deviation 1 For every sample, we do $\frac{(x_i^{(i)} - \mu^{(i)})}{\sigma^{(i)}}$	<h3>Backpropagation</h3> <ul style="list-style-type: none"> We should add an extra feature with a fixed value of 1 to the input. Associate it with weight $b=0$, b bias term
<h3>Standardization vs. Normalization</h3> <ul style="list-style-type: none"> Treat scaling as a hyperparameter and evaluate both normalization and standardization Standardization is generally more robust to outliers than normalization Some guidelines: <ul style="list-style-type: none"> Use standardization for unsupervised learning tasks Use standardization if features are approximately normally distributed Prefer standardization in the presence of outliers Otherwise, use normalization 	<ul style="list-style-type: none"> Standardization is generally more robust to outliers than normalization Some guidelines: <ul style="list-style-type: none"> Use standardization for unsupervised learning tasks Use standardization if features are approximately normally distributed Prefer standardization in the presence of outliers Otherwise, use normalization 	<h3>As the perceptron generates multiple outputs simultaneously, it performs multiple binary predictions, making it a multilabel classifier</h3>
<h3>Handling Missing Values</h3> <ul style="list-style-type: none"> Missing values refer to the absence of data points or entries in a dataset where a value is expected Handling missing values: <ul style="list-style-type: none"> Drop Examples <ul style="list-style-type: none"> Feasible if the dataset is large and outcome is unaffected Drop Features <ul style="list-style-type: none"> Suitable if it does not impact the project's outcome Use algorithms handling missing data Data Imputation <ul style="list-style-type: none"> Replace missing values with computed values 	<ul style="list-style-type: none"> Handling missing values: <ul style="list-style-type: none"> Drop Examples <ul style="list-style-type: none"> Feasible if the dataset is large and outcome is unaffected Drop Features <ul style="list-style-type: none"> Suitable if it does not impact the project's outcome Use algorithms handling missing data Data Imputation <ul style="list-style-type: none"> Replace missing values with computed values 	<h3>Some Notations</h3> <ul style="list-style-type: none"> X is the input data matrix where each row corresponds to an example and each column represents one of the D features W is the weight matrix, structured with one row per input (feature) and one column per neuron Bias terms can be represented separately. Here, b is a vector with a length equal to the number of neurons
<h3>Data Imputation</h3> <ul style="list-style-type: none"> Process of replacing missing values in a dataset with substituted values, typically using statistical or machine learning methods One strategy can be replacing missing values with mean or median of the attribute Cons: ignores feature correlations and complex relationships Mode imputation: replace missing values with the most frequent value Data imputation relies on several assumptions <ul style="list-style-type: none"> Randomness: many methods assume missingness is unrelated to any data Model bias: Incorrect randomness assumptions can lead to biased estimates and flawed conclusions Information Loss: Imputation can obscure patterns, leading to loss of valuable information for advanced models Special Value Method: Replace missing values with a value outside the normal range <ul style="list-style-type: none"> Objective: Enable the learning algorithm to recognize and appropriately handle missing values Middle-Range Imputation: Replace missing values with a value in the middle of the normal range 	<h3>Neural Networks</h3> <ul style="list-style-type: none"> We now investigate a family of machine learning models that draw inspiration from the structure and function of biological neural networks found in animals Computation with neurons 	<h3>Multilayer Perceptron (MLP)</h3> <ul style="list-style-type: none"> Includes an input layer and one or more layers of threshold logic units Layers that are neither input nor output are termed hidden layers
<h3>Feedforward Neural Network (FNN)</h3> <ul style="list-style-type: none"> Information in this architecture flows unidirectionally - from left to right, moving from input to output 	<h3>Activation Functions</h3> <ul style="list-style-type: none"> Digital computations can be broken down into a sequence of logical operations, enabling neural networks to execute any computation Threshold Logic Unit 	<h3>Step 1: Initialization</h3> <ul style="list-style-type: none"> Initialize the weights and biases of the neural network Zero Initialization <ul style="list-style-type: none"> All weights initialized to zero Symmetry problems, all neurons produce identical outputs, preventing effective learning Random Initialization <ul style="list-style-type: none"> Weights are initialized randomly, often using a uniform or normal distribution Breaks the symmetry between neurons, allowing them to learn If not scaled properly, leads to slow convergence or vanishing/exploding gradients
		<h3>Step 2: Forward Pass</h3> <ul style="list-style-type: none"> For each example in the training set Input Layer: Pass input features to first layer Hidden Layer: For each hidden layer, compute the activations (output) by applying the weighted sum of inputs plus bias, followed by an activation function Output Layer: Same process as hidden layers. Output layer activations represent the predicted values
		<h3>Step 3: Compute Loss</h3> <ul style="list-style-type: none"> Calculate the loss (error) using a suitable loss function by comparing the predicted values to the actual target values
		<h3>Step 4: Backward Pass</h3> <ul style="list-style-type: none"> Output Layer: Compute the gradient of the loss with respect to the output layer's weights and biases using the chain rule of calculus Hidden Layers: Propagate the error backward through the network, layer by layer. For each layer, compute the gradient of the loss with respect to the weights and biases. Use the derivative of the activation function to help calculate the gradients Update Weights and Biases: Adjust the weights and biases using the calculated gradients and a learning rate, which determines the step size for each update
		<h3>Activation Function</h3> <ul style="list-style-type: none"> The training algorithm, known as backpropagation, employs gradient descent, necessitating the calculation of the partial derivatives of the loss function The step function in the multilayer perceptron had to be replaced, as it consists only of flat surfaces. Gradient descent cannot progress on flat surfaces due to their zero derivative Nonlinear activation functions are paramount because, without them, multiple layers in the network would only compute a linear function of the inputs

Gradient Descent: An optimization algorithm used to minimize the loss function by iteratively moving towards the steepest descent as defined by the negative of the gradient

Vanishing Gradients

- Gradients become too small, hindering weight updates
- Stalled neural networks research in early 2000s
- Sigmoid and its derivative were key factors
- Common initialization: Weights/biases from $N(0,1)$ contributed to the issue
- Solutions:
 - Alternative activation functions
 - Weight initialization

Output Layer

- For regression tasks:
 - 1 output neuron per dimension
 - Activation function
 - ReLU/softplus if positive
 - Sigmoid/tanh if bounded
 - Loss: mean squared error
- For classification tasks:
 - 1 output neuron if binary, otherwise 1 per class for multiclass
 - Activation function
 - Sigmoid if binary or multi-label
 - Softmax if multi-class
 - Loss: cross-entropy

Softmax

- An activation function used in multi-class classification problems to convert a vector of raw scores into probabilities that sum to 1
- Given a vector $Z = [z_1, z_2, \dots, z_n]$

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

$\sigma(z)_i$: probability of i^{th} class, n is number of classes

Cross-Entropy Loss Function

- For one example

$$J(W) = -\sum_{k=1}^K y_k \log(\hat{y}_k)$$

Where:

- K number of classes
- y_k : true prediction for class k
- \hat{y}_k : predicted probability of class k

For a dataset with N examples, the average cross-entropy loss over all examples is computed as

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

Where:

- i indexes over the different examples in the dataset
- $y_{i,k}$ and $\hat{y}_{i,k}$ are the true and predicted probabilities for class k of example i, respectively

Regularization

Comprises a set of techniques designed to enhance a model's ability to generalize by mitigating overfitting. By discouraging excessive model complexity, these methods improve the model's robustness and performance on unseen data

In numerical optimization, it is standard practice to incorporate additional terms into the objective function to deter undesirable model characteristics

For a minimization problem, the optimization process aims to circumvent the substantial costs associated with these penalty terms

Norm

Assigns a non-negative length to a vector

L_p norm of a vector $Z = [z_1, z_2, \dots, z_n]$ is defined as

$$\|Z\|_p = \left(\sum_{i=1}^n |z_i|^p \right)^{\frac{1}{p}}$$

L_1 norm:

$$\|Z\|_1 = \sum_{i=1}^n |z_i|$$

L_2 norm:

$$\|Z\|_2 = \sqrt{\sum_{i=1}^n z_i^2}$$

L_1 and L_2 Regularization

$$MAE(X, W) = \frac{1}{N} \sum_{i=1}^N \|w_i(x_i) - y_i\|_1 + \alpha \|w_i\|_2$$

α and β determine the degree of regularization applied

L_1 regularization:

- Promotes sparsity, setting many weights to zero
- Useful for feature selection by reducing feature reliance

L_2 regularization:

- Promotes small, distributed weights for stability
- Ideal when all features contribute and reducing complexity is key

Dropout

Regularization techniques in neural networks where randomly selected neurons are ignored during training, reducing overfitting by preventing co-adaptation of features

During each training step, each neuron in a dropout layer has a probability p of being excluded from the computation

While seemingly counterintuitive, this approach prevents the network from depending on specific neurons, promoting the distribution of learned representations across multiple neurons

One of the most popular and effective methods for reducing overfitting

Typical improvement in performance is modest

Early Stopping

A regularization technique that halts training once the model's performance on a validation set begins to degrade, preventing overfitting by stopping before the model learns noise

Hierarchy of Concepts

Each layer detects patterns from the output of the layer preceding it

In other words, proceeding from the input to the output of the network, the network uncovers "patterns of patterns"

Start with one layer, then increase the number of layers until the model starts overfitting the training data

Finetune the model adding regularization

With a feed-forward network, the number of parameters grows rapidly with more layers

Convolutional Neural Network

- Specializes in pattern recognition
- Use filters to recognize patterns
- Crucial pattern information is often local
- Convolutional layers reduce parameters significantly

Unlike dense layers, neurons in a convolutional layer are not fully connected to the preceding layer

Neurons connect only within their receptive fields

Kernel

- A small matrix that slides over the input data to perform convolution
- With the kernel placed over a specific region of the input matrix, the convolution is element-wise multiplication followed by a summation of the results to produce a single scalar value
- In convolutional networks, the kernels are automatically learned by the network

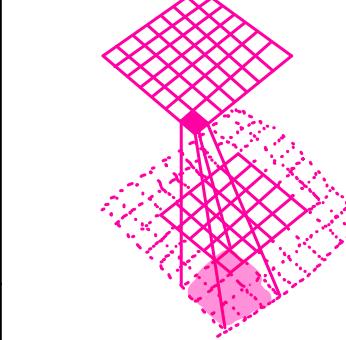
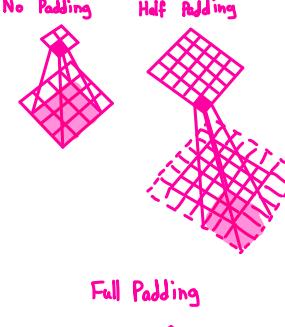
Receptive Field

Each unit is connected to neurons in its receptive fields

Unit i,j in layer L is connected to the units i to $i+f_h-1$, j to $j+f_w-1$ of the layer $L-1$, where f_h and f_w are respectively the height and width of the receptive field.

Padding

Zero padding: In order to have layers of the same size, the grid can be padded with zeros



Stride

It is possible to connect a larger layer ($L-1$) to a smaller one (L) by skipping units. The number of units skipped is called stride, S_h and S_w .

Unit i,j in layer L is connected to the units i to $i+S_h+f_h-1$, j to $j+S_w+f_w-1$ of the layer $L-1$, where f_h and f_w are respectively the height and width of the receptive field. S_h and S_w are respectively the height and weight strides.

Filters

A window of size $f_h \times f_w$ is moved over the output of layers $L-1$, referred to as the input feature map, position by position

For each location, the product is calculated between the extracted patch and a matrix of the same size, known as a convolutional kernel or filter. The sum of the values in the resulting matrix constitutes the output for that location

Model

For every node, we have that its output is

$$Z_{i,j,k} = b_k + \sum_{u=0}^{f_h-1} \sum_{v=0}^{f_w-1} \sum_{l=0}^{S_h-1} \sum_{m=0}^{S_w-1} X_{i+l,f_h+u,j+m,f_w+v} \cdot W_{u,v,l,k}$$

$$i' = i \times S_h + u \quad j' = j \times S_w + v$$

Convolutional Layer

- A layer full of neurons using the same filter outputs a feature map
- During training the convolutional layer will automatically learn the most useful filters for its task
- The layers above will learn to combine them into more complex patterns
- The fact that all neurons in a feature map share the same parameters dramatically reduces the number of parameters in the model

Summary of CNNs

- Feature Map:** In convolutional neural networks, the output of a convolution operation is known as a feature map. It captures the features of the input data as processed by a specific kernel
- Kernel Parameters:** The parameters of the kernels are learned through the backpropagation process, allowing the network to optimize its feature extraction capabilities based on the training data
- Bias Term:** A single bias term is added uniformly to all entries of the feature map. This bias helps adjust the activation level, providing additional flexibility for the network to better fit the data
- Activation Function:** Following the addition of the bias, the feature map values are typically passed through an activation function

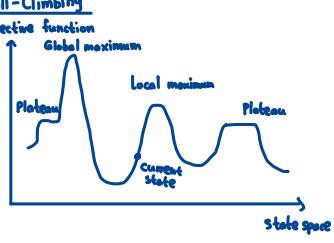
Pooling

- A pooling layer exhibits similarities to a convolutional layer
- Each neuron in a pooling layer is connected to a set of neurons within a receptive field
- However, unlike convolutional layers, pooling layers do not possess weights
- Instead, they produce an output by applying an aggregating function, commonly max or mean
- This subsampling process leads to a reduction in network size; each window of dimensions $f_h \times f_w$ is condensed to a single value, typically the maximum or mean of that window
- Dimensionality Reduction: Pooling layers reduce the spatial dimensions (width and height) of the input feature maps. This reduction decreases the number of parameters and computational load in the network, which can help prevent overfitting
- Feature Extraction: By summarizing the presence of features in a region, pooling layers help retain the most critical information while discarding less important details. This process enables the network to focus on the most salient features
- Translation Invariance: Pooling introduces a degree of invariance to translations and distortions in the input
- Noise Reduction: Pooling can help smooth out noise in the input by aggregating information over a region, thus emphasizing consistent features over random variations
- Hierarchical Feature Learning: By reducing the spatial dimensions progressively through the layers, pooling layers allow the network to build a hierarchical representation of the input data, capturing increasingly abstract and complex features at deeper layers

Search

When the correct action to take is not immediately obvious, an agent may need to plan ahead: to consider a sequence of actions that form a path to a goal state. Such an agent is called a problem-solving agent, and the computational process it undertakes is called search

An agent is an entity that performs actions. A "normal agent" is one that acts to achieve the "best" outcome. Conceptually, an agent perceives its environment through sensors and interacts with it using actuators

<h3>Search Environment Characteristics</h3> <ul style="list-style-type: none"> Observability: Partially observable, or fully observable Agent Composition: Single or multiple agents Predictability: Deterministic or non-deterministic State Dependency: Stateless or stateful Temporal Dynamics: Static or Dynamic State Representation: Discrete or Continuous 	<h3>Depth-First Search</h3> <ul style="list-style-type: none"> Uses a stack to keep track of states The Algorithm <p>Input: Initial and Goal States Initialize a stack that stores the frontier states Initialize a set that stores the explored nodes while the stack is not empty: Top state of stack is current state If the current state is the goal state then return the path Otherwise, add all its neighbors to the explored set and the stack return None if no solution found</p>	<h3>Measuring Performance</h3> <ul style="list-style-type: none"> Completeness: Does the algorithm ensure that a solution will be found if one exists, and accurately indicate failure when no solution exists? Cost Optimality: Does the algorithm identify the solution with the lowest path count among all possible solutions? Time Complexity: How does the time required by the algorithm scale with respect to the number of states and actions? Space Complexity: How does the space required by the algorithm scale with respect to the number of states and actions?
<h3>Search Problem Definition</h3> <ul style="list-style-type: none"> A collection of states, referred to as the state space An initial state where the agent begins One or more goal states that define successful outcomes A set of actions available in a given state s A transition model that determines the next state based on the current state and selected action An action cost function that specifies the cost of performing action a in state s to reach state s' 		
<h3>Paths and Solutions</h3> <ul style="list-style-type: none"> A path is defined as a sequence of actions A solution is a path that connects the initial state to the goal state An optimal solution is the path with the lowest cost among all possible solutions We assume that the path cost is the sum of the individual action costs, and all costs are positive. The state space can be conceptualized as a graph, where the nodes represent the states and the edges correspond to the actions 	<h3>Heuristic Search</h3> <ul style="list-style-type: none"> Use domain-specific knowledge regarding the goal's state location Let $f(n)$ be a heuristic function that estimates the cost of the cheapest path from the current state or node n to the goal In most finding problems, one might employ the straight-line distance from the current node to the destination as heuristic. Although an actual path may not exist along that straight line, the algorithm will prioritize expanding the node closest to the destination based on this straight-line measurement This approach is known as best-first search To implement this, we can employ a priority queue, which is sorted according to the values of the heuristic function $h(n)$ 	<h3>Local Search</h3> <ul style="list-style-type: none"> Local search algorithms operate by searching from a start state to neighboring states, without keeping track of the paths, nor the set of states that have been reached Optimizes memory utilization while effectively solving problems in extensive or infinite state spaces Find the "best" state according to an objective function, thereby locating the global maximum Commonly referred to as hill climbing
<h3>Search Tree</h3> <ul style="list-style-type: none"> A conceptual tree structure where nodes represent states in a state space, and edges represent possible actions, facilitating systematic exploration to find a path from an initial state to a goal state The root of the search tree represents the initial state of the problem Expanding a node involves evaluating all possible actions available from that state The result of an action is the new state achieved after applying that action to the current state Similar to other tree structures, each node (except root and leaf nodes) has a parent and may have children Any state corresponding to a node in the search tree is considered reached Frontier nodes are those that have been reached but have not yet been expanded 	<h3>A* Search</h3> <ul style="list-style-type: none"> The most common informed search $f(n) = g(n) + h(n)$ where $g(n)$ is the path cost from the initial state to n $h(n)$ is an estimate of the cost of the shortest path from n to the goal state A* guarantees a solution if one exists Cost optimality requires an admissible heuristic An admissible heuristic never overestimates the goal cost 	<h3>Hill-Climbing</h3>  <p>Given as input a problem</p> <ul style="list-style-type: none"> Current is the initial state of problem while not done do <ul style="list-style-type: none"> neighbour is the highest-valued successor state of current if $\text{value(neighbour)} \leq \text{value(current)}$ then return current set current to neighbour
<h3>Uninformed Search</h3> <ul style="list-style-type: none"> A search strategy that explores the search space using only the information available in the problem definition, without any domain-specific knowledge, evaluating nodes based solely on their inherent properties rather than estimated costs or heuristics 	<h3>Best-First Search</h3> <ul style="list-style-type: none"> The Algorithm: <p>First initialize a priority queue frontier Then we initialize the heuristic to the distance between start and end state Push the initial state with the heuristic value into the queue</p> <p>Create a set that stores explored states</p> <p>While the queue is not empty if the current state is the goal state, then return the successful path Otherwise, we add the current state to the explored set</p> <p>We expand the current state and go through each neighbor</p> <p>For every unexplored neighbor, we increment the path cost, calculate the cost to the goal, update total cost, then add it to the priority queue and the set of explored states</p>	<h3>8-Queens Problem</h3> <ul style="list-style-type: none"> Involves placing eight queens on an 8x8 chessboard such that no two queens threaten each other, meaning no two queens share the same row, column, or diagonal A grid representation of the current state permits the illegal placement of two queens in the same column Instead, we can represent the state as a list (state), where each element corresponds to the row position of the queen in its respective column $\text{state}[i]$ is the row of the queen in column i
<h3>Breadth-First Search</h3> <ul style="list-style-type: none"> Employs a queue to manage the frontier nodes, which are also known as the open list It finds the shortest path from the initial state to the goal state The Algorithm: <p>Input: Initial and goal states we initialize a queue that stores the frontier nodes we also initialize a set that stores explored states Add the initial state to both While the queue is not empty</p> <p>Current state and path is first element of queue if the current state is the goal state, then return the successful path Otherwise, we expand the current state and we add the neighbor states to the set of explored states and the queue return None if no solution found</p>	<p>State = [0, 4, 7, 5, 2, 6, 1, 3]</p> <ul style="list-style-type: none"> We can randomly create an initial state Representations <ul style="list-style-type: none"> Unconstrained Placement: $\binom{64}{8} = 4,426,165,318$ possible configurations Column Constraint: Uses a list of length 8, with each entry indicating the row of a queen in its respective column, resulting in 8! = 40,320 configurations Row and Column Constraints: Model board states as permutations of the 8 row indices, reducing configurations to 8! = 40,320 	<h3>Travelling Salesman Problem</h3> <p>A classic optimization problem that seeks the shortest possible route for a salesman to visit a set of cities, returning to the origin city, while visiting each city exactly once</p> <ul style="list-style-type: none"> Use a list where each element represents the index of a city, and the order of elements indicates the sequence of city visits Ways to generating a neighboring solution <ul style="list-style-type: none"> Swap Two Cities: Select two cities at random and swap their positions <ul style="list-style-type: none"> Pros: Simple and effective for exploring nearby solutions Cons: Change may be too small, potentially slowing down convergence Reverse Segment: Select two indices and reverse the segment between them <ul style="list-style-type: none"> Pros: More effective at finding shorter paths compared to simple swaps Cons: Can still be computationally expensive as the number of cities increases Remove and Reconnect: Removes three edges from the route and reconnects the segments in the best possible way. This can generate up to 7 different routes

- Pros: Provides more extensive changes and can escape local optima more effectively than 2-opt
- Cons: More complex and computationally expensive to implement
- Insertion Move: Select a city and move it to a different position in the route
 - Pros: Offers a balance between small and large changes, making it useful for fine-tuning solutions
 - Cons: May require more iterations to converge to an optimal solution
- Shuffle Subset: Select a subset of cities in the route and randomly shuffle their order
 - Pros: Introduces larger changes and can help escape local minima
 - Cons: Can lead to less efficient routes if not handled carefully

Simulated Annealing and TSP

- We can use simulated annealing to solve the travelling salesman problem
- The inputs are: a distance matrix containing the distances between cities, the initial temperature, the cooling rate, and maximum number of iterations
- We first determine the number of cities, the current route, and the current cost
- We initialize the best route and best cost to the current route and current cost
- We initialize the temperature to the initial temperature
- We iterate for a while. For each iteration, we do
 - Pick a neighbor and its cost
 - If the neighbor costs less, then we move to it. Otherwise, we move to it with a probability
 - We then cool down the temperature
- As $t \rightarrow \infty$, this algorithm exhibits behaviour characteristic of a random walk. During this phase, any neighboring state, regardless of whether it improves the objective function, is accepted. This facilitates exploration and occurs at the start of the algorithm's execution
- As $t \rightarrow 0$, the algorithm behaves like hill climbing. In this phase, only those states that enhance the objective function's value are accepted, ensuring that the algorithm consistently moves towards optimal solutions - specifically, towards lower values in minimization problems. This phase emphasizes the exploitation of promising solutions and occurs towards the algorithm's conclusion