

# CSI 4106

## Artificial Intelligence

### Study Guide



### Fall 2024

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

- Predictions are made as follows:

$$\begin{aligned} \cdot y_i = 0, & \text{ if } h_\theta(x_i) < 0.5 \\ \cdot y_i = 1, & \text{ if } h_\theta(x_i) \geq 0.5 \end{aligned}$$

- The values of  $\theta$  are learned using gradient descent

### Feature Engineering

- The process of creating, transforming, and selecting variables from raw data to improve the performance of machine learning models

### Underfitting and Overfitting

#### Underfitting

- Model is too simple
- Uninformative features
- Poor performance on both train and test data

#### Overfitting

- Model is too complex
- Too many features
- Excellent performance on training set, but poor performance on test set

### Learning Curves

- 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

### Bias / Variance Tradeoff

- 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

### Confusion Matrix

		Positive (predicted)	Negative (predicted)
Positive (Actual)	True Positive (TP)	False Negative (FN)	
	False Positive (FP)	True Negative (TN)	

### Accuracy

- The accuracy of a model measures how accurate the result is

- Ratio of correct number of predictions to total number of predictions

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

- 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

### Precision

- Also known as positive predictive value
- Proportion of true positive predictions among all positive predictions

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

- An algorithm that makes a small number of high-confidence predictions might achieve a high precision score, but may not be helpful

### Recall

- Also known as sensitivity or true positive rate (TPR)

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

- Proportion of true positive instances correctly identified among all actual positive instances

### F<sub>1</sub> Score

- Harmonic mean of precision and recall

$$\begin{aligned} F_1 \text{ score} &= \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}} \\ &= 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \\ &= \frac{TP}{\frac{TP+FP+FN}{2}} \end{aligned}$$

### Micro Performance Metrics

- Micro performance metrics aggregate the contributions of all classes to compute the average performance metric, such as precision, recall, or  $F_1$  score.

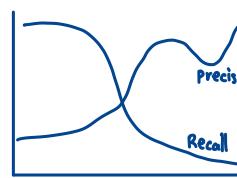
- This approach treats each individual prediction equally, providing a balanced evaluation by emphasizing the performance on frequent classes

### Macro Performance Metrics

- Compute the performance metric independently for each class and then average these metrics.
- Treats each class equally, regardless of frequency

### Precision-Recall Tradeoff

- 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 decreases, fewer examples are classified as positive, which may result in the classifier predicting no positive instances at all



Threshold

### ROC Curve

Receiver Operating Characteristics curve

- 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

$$\text{TPR} = \frac{TP}{TN+FP}$$

- TPR approaches one when the number of false negative predictions is low

$$\text{FPR} = \frac{FP}{FP+TN}$$

- FPR approaches zero when the number of false positive predictions is low