

# CSC413 Neural Networks and Deep Learning

## Lecture 1: Introduction and Review of Linear Models

January 9 / 11, 2024

# Table of Contents

- 1 Welcome to CSC413!
- 2 How can we make an intelligent agent?
- 3 Overview of Supervised Learning with Linear Models
- 4 Deep Learning
- 5 Logistics
- 6 What to do this week?

# Section 1

Welcome to CSC413!

## **Amanjit Singh Kainth**

- LEC 0101/2001: Tuesday (T) 1:00PM - 4:00PM
- Office Hours: TBA

## **Amir-massoud Farahmand**

- LEC 0201/2101: Thursday (R) 1:00PM - 4:00PM
- Office Hours: TBA

## **Robert (Rupert) Wu**

- LEC 5101/2501: Tuesday (T) 6:00PM - 9:00PM
- Office Hours: T4-5 BA 2272

## Section 2

How can we make an intelligent agent?

# How can we make an intelligent agent?

- What does it mean to have an intelligent agent?
- What do we need to create it?



Figure 1: An agent ...

# An AI Agent



Figure 2: ... observes the world ...



# An AI Agent



Figure 3: ... takes an action and its states changes ...

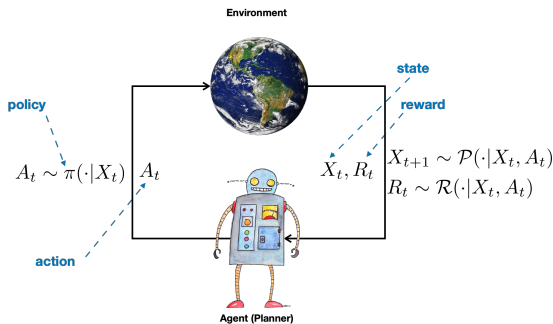


Figure 4: ... with the goal of achieving long-term rewards.

# An AI Agent: Some Requirements

This agent has to

- **predict** how the world works, e.g.,
  - classify different objects
  - estimate the probability of certain events happening in the future
- **plan** its actions in order to achieve its long-term goals



We use Machine Learning and Neural Networks to move towards this goal.

# What is the difference between. . .

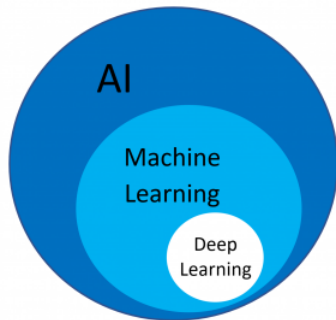
- Artificial Intelligence
- Machine Learning
- Deep Learning

# AI vs ML vs DL

**Artificial Intelligence:** Create intelligent machines that perceive, reason, and act like humans. (CSC384)

**Machine Learning:** Design algorithms to automatically learn from data. (CSC311)

**Deep Learning:** Using deep neural networks to automatically learn from data. (CSC413)



# Why machine learning after all?

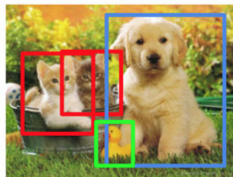
Our agent lives in a complicated world. It is difficult to program the correct behaviour (ex. recognizing chairs, apples, etc.) manually.

**Machine learning approach:** Write an algorithm to automatically learn from data.



Gary Chavez added a photo you might ...  
be in.

about a minute ago · 👤



CAT, DOG, DUCK

(ALREADY FAMILIAR EXAMPLE)

I'M GOING TO THE THEATER = ICH GEHE INS THEATER

I'M GOING TO THE CINEMA = ICH GEHE INS KINO

→ KINO ←

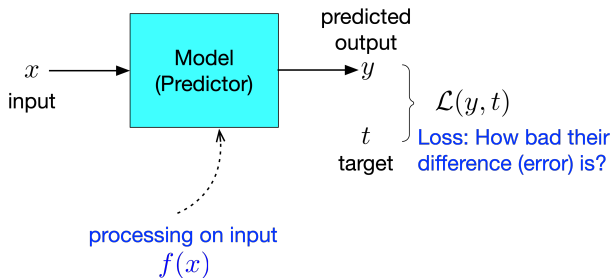
# Types of problems in Machine Learning

- **Supervised Learning:** have labeled examples of the correct behaviour, i.e. ground truth input/output response
  - Regression (e.g. height prediction)
  - Classification (e.g. sentiment classification)
- **Unsupervised Learning:** no labeled examples; instead, looking for interesting patterns in the data (e.g. clustering)
- **Reinforcement Learning:** learning system receives a reward signal, tries to learn to maximize the reward signal (e.g. playing StarCraft)

This categorization is not strict. There are overlaps between these problems, and there are problems that belong to more than one category.

- Examples: semi-supervised learning, self-supervised learning, model learning for reinforcement learning, etc.

# The Machine Learning Approach



Reframe **learning problems** into **optimization** problems by:

- Choosing a model (with parameters to be tuned)
- Choosing a loss/cost function to measure how well the model fits the data given a choice of parameters
- Choosing an optimizer to minimize the cost function

Different machine learning approaches differ in the model, loss, and optimizer choice.



## Section 3

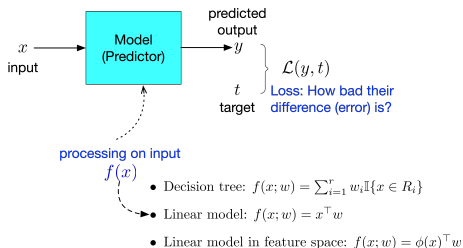
# Overview of Supervised Learning with Linear Models

# Overview of Supervised Learning with Linear Models

This is going to be quick! We assume that you are already familiar with these material because you have taken CSC311.

# Supervised Learning with Linear Models

- Let us review Linear Models.
- They also appear as one of the building blocks in neural networks and deep learning.



# Common supervised learning problems

Recall the types of supervised learning problems:

- **Regression**: predict a scalar-valued target (e.g. stock price)
- **Classification**: predict a label
  - **Binary classification**: predict a binary label (e.g. spam vs. non-spam email)
  - **Multi-class classification**: predict a discrete label (e.g. object category, from a list)

# Problem Setup: Regression

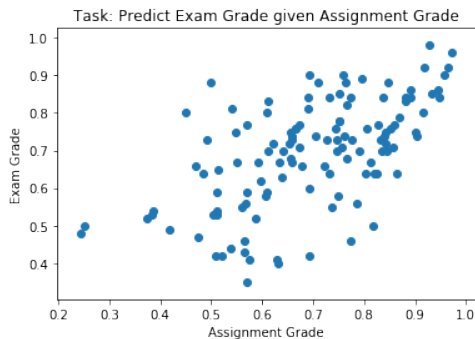
**Input:** Represented using the vector  $\mathbf{x}$

- Example:  $\mathbf{x}$  represents assignment grades (0-100)
- To start, let's assume that  $\mathbf{x}$  is a scalar, and that we only have the cumulative assignment grade

**Output:** Represented using the scalar  $t$

- Example:  $t$  represents the grade on an exam (0-100)
- We'll use the scalar  $y$  to denote a *prediction* of the value of  $t$

# Example: Exam Grade Prediction



- Data:  $(\mathbf{x}^{(1)}, t^{(1)})$ ,  $(\mathbf{x}^{(2)}, t^{(2)})$ ,  $\dots$   $(\mathbf{x}^{(N)}, t^{(N)})$
- The  $\mathbf{x}^{(i)}$  are called *inputs*
- The  $t^{(i)}$  are called *targets*

# Regression with a Linear Model

- A **model** implicitly or explicitly encodes our assumptions about the underlying nature of the data we wish to learn about.
- But recall the adage that *All models are wrong, but some are useful.*
  - In ML, we often choose models without really explicitly thinking about our assumptions about the data generation process.
- The **model**, or **architecture** defines the set of allowed family of **hypotheses**.

In linear regression, our **model** looks like

$$y = \sum_j w_j x_j + b,$$

where  $y$  is a prediction for  $t$ , and the  $w_j$  and  $b$  are **parameters** of the model, to be determined based on the data.

# Linear Regression for Exam Grade Prediction

For the exam prediction problem, we only have a single feature, so we can simplify our model to:

$$y = wx + b$$

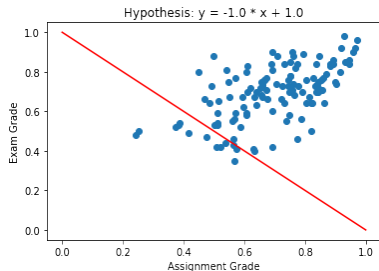
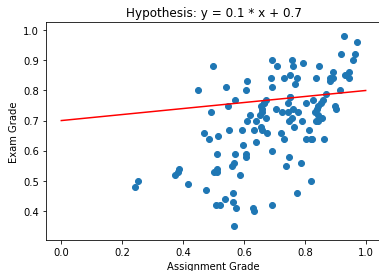
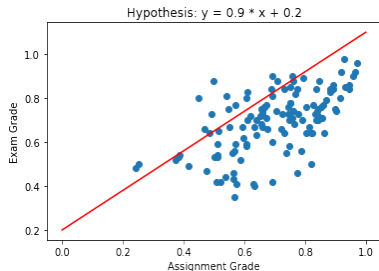
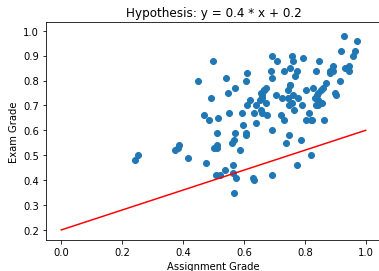
Our **hypothesis space** includes all functions of the form  $y = wx + b$ . Here are some examples:

- $y = 0.4x + 0.2$
- $y = 0.9x + 0.2$
- $y = 0.1x + 0.7$
- $y = -x - 1$
- ...

The variables  $w$  and  $b$  are called **weights** or **parameters** of our model. (Sometimes  $w$  and  $b$  are referred to as coefficients and intercept/bias, respectively.)

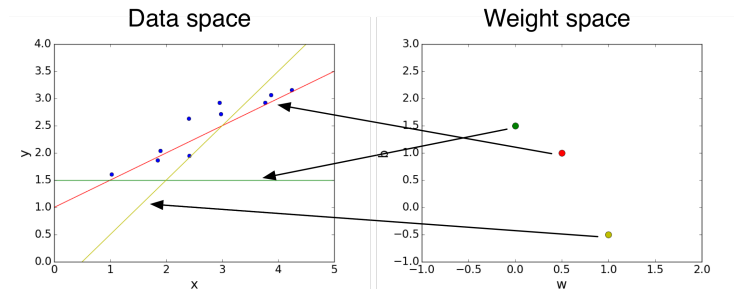


# Which hypothesis is better suited to the data?



# Hypothesis Space

We can visualize the hypothesis space or **weight space**:



Each *point* in the weight space represents a hypothesis.

# Cost Function (Loss Function)

The “badness” of an entire hypothesis is the average badness across our labeled data.

$$\begin{aligned}\mathcal{E}(w, b) &= \frac{1}{N} \sum_i \mathcal{L}(y^{(i)}, t^{(i)}) \\ &= \frac{1}{2N} \sum_i (y^{(i)} - t^{(i)})^2 \\ &= \frac{1}{2N} \sum_i ((wx^{(i)} + b) - t^{(i)})^2\end{aligned}$$

This is called the **cost** of a particular hypothesis (in practice, “loss” and “cost” functions are used inter-changeably).

Since the loss depends on the choice of  $w$  and  $b$ , we call  $\mathcal{E}(w, b)$  the **cost function**.

# Minimize Cost: Direct Solution

Find a *critical point* by setting

$$\frac{\partial \mathcal{E}}{\partial \mathbf{w}} = 0$$
$$\frac{\partial \mathcal{E}}{\partial \mathbf{b}} = 0$$

Possible for our hypothesis space, and covered in the notes.

However, let's use a technique that can also be applied to more general models.

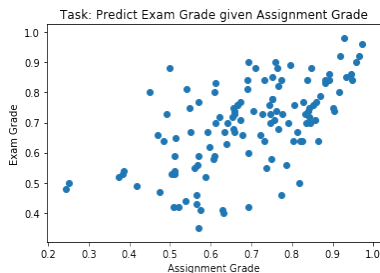
# Minimize Cost: Gradient Descent

We can use gradient descent to minimize the cost function.

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \mathcal{E}}{\partial \mathbf{w}}$$
$$\frac{\partial \mathcal{E}}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial \mathcal{E}}{\partial w_1} \\ \dots \\ \frac{\partial \mathcal{E}}{\partial w_D} \end{bmatrix}$$

The  $\alpha$  is the **learning rate**, which we choose.

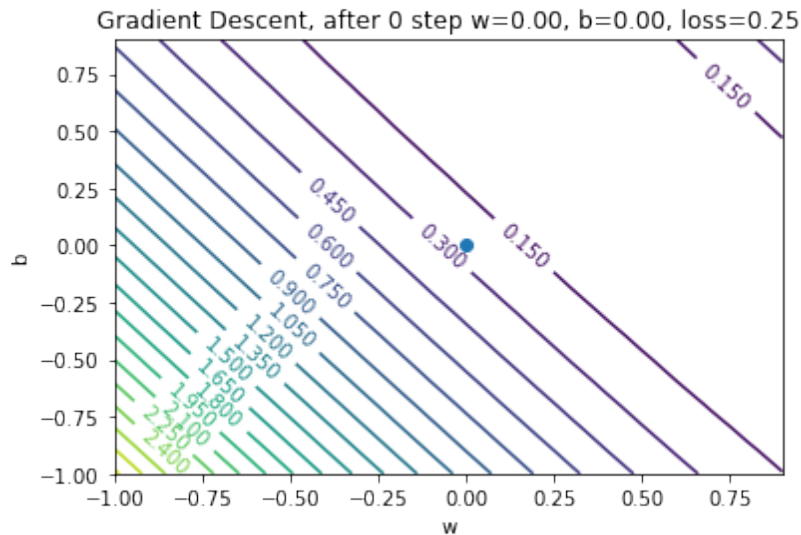
# Gradient Descent for Grade Prediction



We'll initialize  $w = 0$  and  $b = 0$  (arbitrary choice)

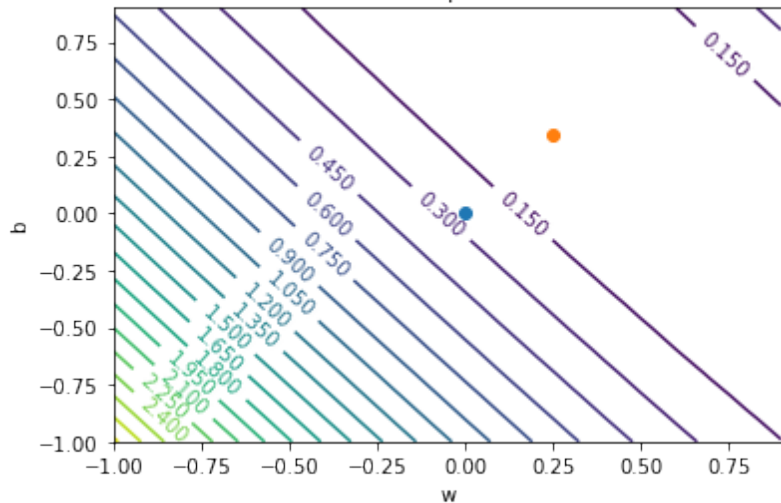
We'll also choose  $\alpha = 0.5$

# Gradient Descent: Step 0



# Gradient Descent: Step 1

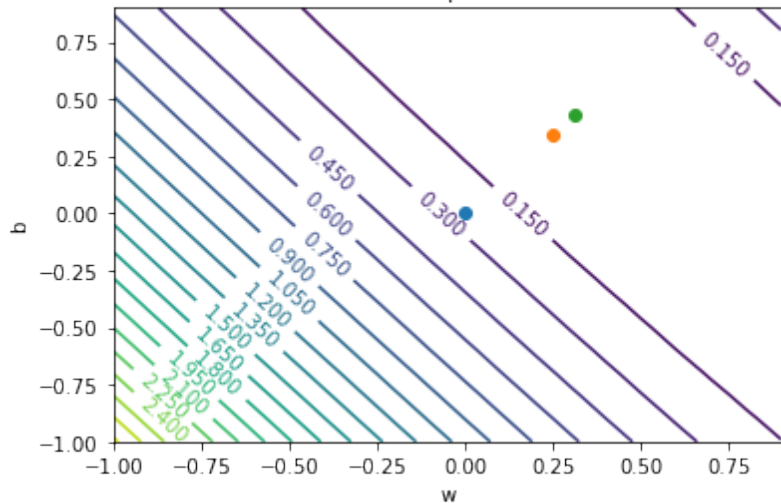
Gradient Descent, after 1 step  $w=0.25$ ,  $b=0.35$ , loss=0.02





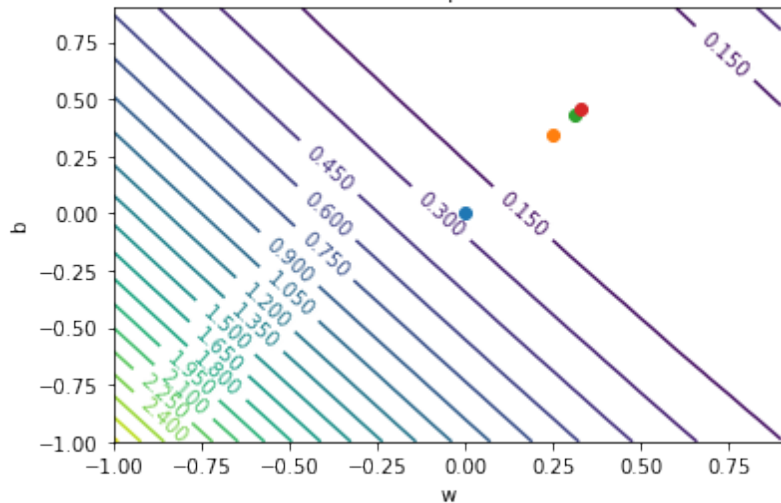
# Gradient Descent: Step 2

Gradient Descent, after 2 step  $w=0.31$ ,  $b=0.43$ ,  $\text{loss}=0.01$



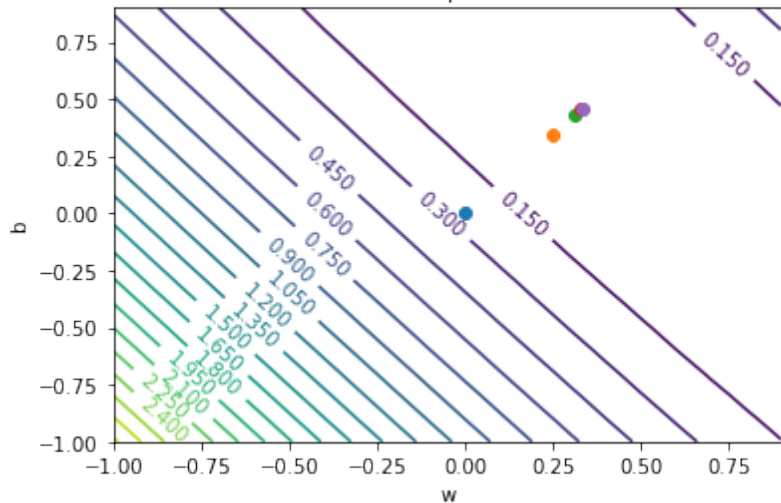
# Gradient Descent: Step 3

Gradient Descent, after 3 step  $w=0.33$ ,  $b=0.45$ , loss=0.01



# Gradient Descent: Step 4

Gradient Descent, after 4 step  $w=0.33$ ,  $b=0.46$ ,  $\text{loss}=0.01$



# Gradient Descent: When to Stop?

In theory:

- Stop when  $w$  and  $b$  stop changing (convergence)

In practice:

- Stop when  $\mathcal{E}$  “almost” stops changing (approximate convergence)
- Stop until we run out of our computational budget or get tired of waiting more

# Gradient Descent: How to Choose the Learning Rate?

- If  $\alpha$  is too small, then training will be *slow*
  - Take a long time to (approximately) converge
- If  $\alpha$  is too large, then we can have divergence!
  - It does not converge

# Computing the Gradient

To compute the gradient  $\frac{\partial \mathcal{E}}{\partial w}$

$$\frac{\partial \mathcal{E}}{\partial w} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}(y^{(i)}, t^{(i)})}{\partial w}$$

But this computation can be expensive if  $N$  is large!

# Computing the Gradient

To compute the gradient  $\frac{\partial \mathcal{E}}{\partial w}$

$$\frac{\partial \mathcal{E}}{\partial w} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}(y^{(i)}, t^{(i)})}{\partial w}$$

But this computation can be expensive if  $N$  is large!

Solution: estimate  $\frac{\partial \mathcal{E}}{\partial w}$  using a *subset* of the data

# Stochastic Gradient Descent

Full batch gradient descent:

$$\frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}(y^{(i)}, t^{(i)})}{\partial w}$$

Stochastic Gradient Descent:

Estimate the above quantity by computing the average of  $\frac{\partial \mathcal{L}(y^{(i)}, t^{(i)})}{\partial w}$  across a small number of  $i$ 's

The set of examples that we use to estimate the gradient is called a **mini-batch**.

The number of examples in each mini-batch is called the **mini-batch size** or just the **batch size**



# Stochastic Gradient Descent Algorithm

In theory, any way of sampling a mini-batch is okay.

In practice, SGD is almost always implemented like this:

```
# repeat until convergence:  
  # randomly split the data into mini-batches of size k  
  # for each mini-batch:  
    # estimate the gradient using the mini-batch  
    # update the parameters based on the estimate
```

- Each pass of the inner loop is called an **iteration**.
  - One iteration = one update for each weight
- Each pass of the outer loop is called an **epoch**.
  - One epoch = one pass over the data set

# Iterations, Epochs, and Batch Size

Suppose we have 1000 examples in our training set.

- Q: How many iterations are in one epoch if our batch size is 10?
- Q: How many iterations are in one epoch if our batch size is 50?

# Batch size choice

- Q: What happens if the batch size is **too large**?
- Q: What happens if the batch size is **too small**?

# Linear Regression Summary

---

Model  $y = \mathbf{w}^\top \mathbf{x} + b$

Loss  $\mathcal{L}(y, t) = (y - t)^2$

Func-  
tion

Optimization  $\min_{\mathbf{w}, b} \mathcal{E}(\mathbf{w}, b)$  via Gradient Descent

Method

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \mathcal{E}}{\partial \mathbf{w}}, b \leftarrow b - \alpha \frac{\partial \mathcal{E}}{\partial b}$$

---

# Linear Regression Vectorization

Use vectors rather than writing

$$\mathcal{E}(\mathbf{w}, b) = \frac{1}{2N} \sum_i ((\mathbf{w}\mathbf{x}^{(i)} + b) - t^{(i)})^2$$

So we have:

$\mathbf{y} = \mathbf{X}\mathbf{w} + b\mathbf{1}$ , where

$$\mathbf{X} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_D^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_D^{(2)} \\ \dots & \dots & \dots & \dots \\ x_1^{(N)} & x_2^{(N)} & \dots & x_D^{(N)} \end{bmatrix}, \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_D \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \dots \\ y^{(N)} \end{bmatrix}, \mathbf{t} = \begin{bmatrix} t^{(1)} \\ t^{(2)} \\ \dots \\ t^{(N)} \end{bmatrix}$$

(You can also fold the bias  $b$  into the weight  $\mathbf{w}$ , but we won't.)

# Vectorized Loss Function

After vectorization, the loss function becomes:

$$\mathcal{E}(\mathbf{w}) = \frac{1}{2N}(\mathbf{y} - \mathbf{t})^\top (\mathbf{y} - \mathbf{t})$$

or

$$\mathcal{E}(\mathbf{w}) = \frac{1}{2N}(\mathbf{X}\mathbf{w} + b\mathbf{1} - \mathbf{t})^\top (\mathbf{X}\mathbf{w} + b\mathbf{1} - \mathbf{t})$$

# Vectorized Gradient Descent

$$b \leftarrow b - \alpha \frac{\partial \mathcal{E}}{\partial b}$$
$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \mathcal{E}}{\partial \mathbf{w}}$$

Where  $\frac{\partial \mathcal{E}}{\partial \mathbf{w}}$  is the vector of partial derivatives:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial \mathcal{E}}{\partial w_1} \\ \dots \\ \frac{\partial \mathcal{E}}{\partial w_D} \end{bmatrix}$$

# Why vectorize?

Vectorization is *not* just for mathematical elegance.

When using Python with numpy/PyTorch/Tensorflow/JAX, code that performs vector computations is faster than code that loops.

Same holds for many other high level languages and software.



# Problem Setup: Classification

- Data:  $(x^{(1)}, t^{(1)}), (x^{(2)}, t^{(2)}), \dots (x^{(N)}, t^{(N)})$
- The  $x^{(i)}$  are called *inputs*
- The  $t^{(i)}$  are called *targets*

In classification, the  $t^{(i)}$  are discrete.

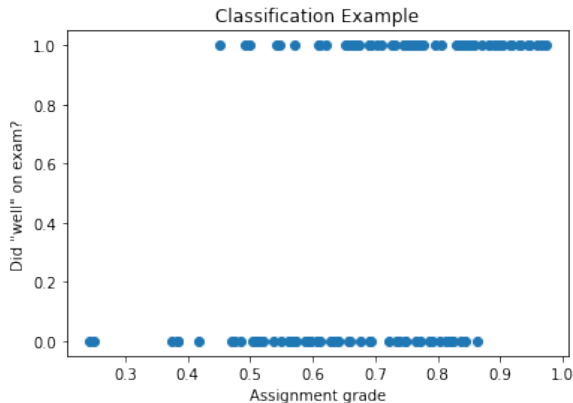
In binary classification, we'll use the labels  $t \in \{0, 1\}$  (or  $\{-1, +1\}$ ).

Training examples with

- $t = 1$  is called a **positive example**
- $t = 0$  is called a **negative example**

# Classification Example

- $x^{(i)}$  represents a person's assignment grade
- $t^{(i)}$  represents whether that person had a "high" exam grade (arbitrary cutoff)



## Q: Why not use regression?

Why can't we set up this problem as a regression problem?

Use the model:

$$y = wx + b$$

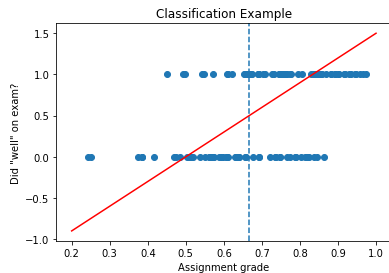
Our prediction for  $t$  would be 1 if  $y \geq 0.5$ , and 0 otherwise.

With the loss function

$$\mathcal{L}(y, t) = \frac{1}{2}(y - t)^2$$

And minimize the cost function via gradient descent?

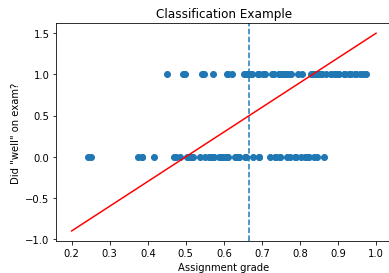
# Classification as Regression: Problem



If we have  $\mathcal{L}(y, t) = \frac{1}{2}(y - t)^2$ , then **points that are correctly classified will still have high loss!**

(blue dotted line above = decision boundary)

# The Problem (continued)



Example: a point on the top right

- Model makes the correct prediction for point on top right
- However,  $(y - t)^2$  is large
- So we are penalizing our model, even though it is making the right prediction!

## Q: Why not use classification error?

Why not still use the model:

$$y = \begin{cases} 1 & \text{if } \mathbf{w}^\top \mathbf{x} + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

But use this loss function instead:

$$\mathcal{L}(y, t) = \begin{cases} 0 & \text{if } y = t \\ 1 & \text{otherwise} \end{cases}$$

## First attempt at a loss function: the 0-1 loss

$$\mathcal{L}(y, t) = \begin{cases} 0 & \text{if } y = t \\ 1 & \text{otherwise} \end{cases}$$

The gradient of this function is 0 almost everywhere!

So gradient descent will not change the weights! We need to define a surrogate loss function that is better behaved.

# Logistic Model

Apply a **non-linearity** or **activation function** to the linear model  $z$ :

$$z = wx + b \quad \text{also called the logit}$$

$$y = \sigma(z) \quad \text{also called a log-linear model}$$

where

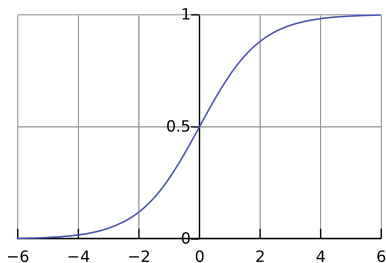
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

is called the logistic or sigmoid function.



# The Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

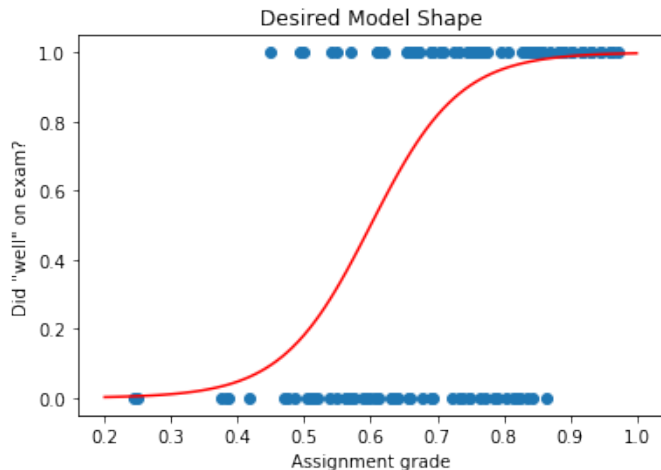


Properties:

- $\sigma(z)$  is between 0 and 1
- $\sigma(0)$  is 0.5

# Logistic Regression Example

A logistic model has this shape:



But how do we train this model?

# Logistic Model with Squared Error Loss?

Suppose we define the model like this:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L}_{SE}(y, t) = \frac{1}{2}(y - t)^2$$

The gradient of  $\mathcal{L}$  with respect to  $w$  is (homework):

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w} &= \frac{\partial \mathcal{L}}{\partial y} \frac{dy}{dz} \frac{\partial z}{\partial w} \\ &= (y - t)y(1 - y)x\end{aligned}$$

# The Problem of Squared Error Loss

Suppose we have a positive example ( $t = 1$ ) that our model classifies extremely wrongly ( $z = -5$ ):

Then we have  $y = \sigma(z) \approx 0.0067$

# The Problem of Squared Error Loss

Suppose we have a positive example ( $t = 1$ ) that our model classifies extremely wrongly ( $z = -5$ ):

Then we have  $y = \sigma(z) \approx 0.0067$

Ideally, the *gradient* should give us strong signals regarding how to update  $w$  to do better.

But...  $\frac{\partial \mathcal{L}}{\partial w} = (y - t)y(1 - y)x$  is small!

# The Problem of Squared Error Loss

Suppose we have a positive example ( $t = 1$ ) that our model classifies extremely wrongly ( $z = -5$ ):

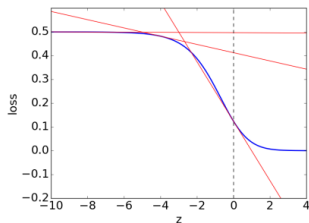
Then we have  $y = \sigma(z) \approx 0.0067$

Ideally, the *gradient* should give us strong signals regarding how to update  $w$  to do better.

But...  $\frac{\partial \mathcal{L}}{\partial w} = (y - t)y(1 - y)x$  is small!

Which means that the update  $w \leftarrow w - \alpha \frac{\partial \mathcal{L}}{\partial w}$  won't change  $w$  much!

# Gradient Signal



The problem with using sigmoid activation with square loss is that we get **poor gradient signal**.

- The loss for a *very wrong prediction* ( $y = 0.0001$ ) vs a wrong prediction ( $y = 0.01$ ) are similar
- This is a problem, because the gradients in the region would be close to 0

We need a loss function that distinguishes between a wrong prediction and a *very wrong prediction*.

# The Cross Entropy Loss

The **cross entropy loss** provides the desired behaviour:

$$\mathcal{L}(y, t) = \begin{cases} -\log(y) & \text{if } t = 1 \\ -\log(1 - y) & \text{if } t = 0 \end{cases}$$

We can write the loss as:

$$\mathcal{L}(y, t) = -t \log(y) - (1 - t) \log(1 - y)$$

We use the Logistic Model with cross entropy loss, the resulting model is called **Logistic Regression** model. Note that it is a classification method and not a regression one, as we use it in this course.



# Logistic Regression Summary

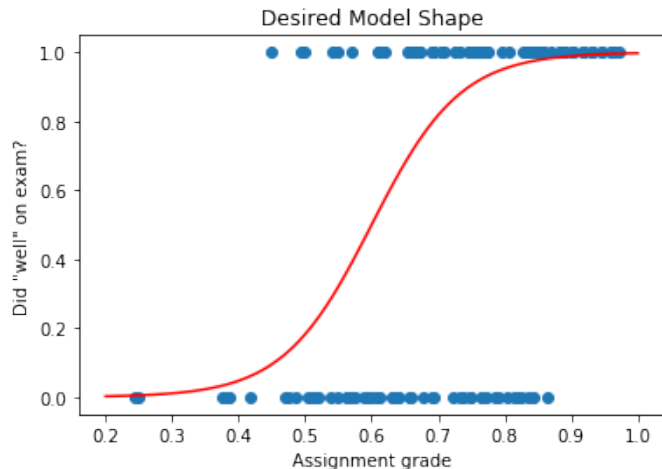
---

Model	$y = \sigma(\mathbf{w}^\top \mathbf{x} + b)$
Loss	$\mathcal{L}(y, t) = -t \log(y) - (1 - t) \log(1 - y)$
Func- tion	
Optimization	$\min_{\mathbf{w}, b} \mathcal{E}(\mathbf{w}, b)$ via Gradient Descent
Method	$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \mathcal{E}}{\partial \mathbf{w}}, b \leftarrow b - \alpha \frac{\partial \mathcal{E}}{\partial b}$

---

# Grade Classification Example

After running gradient descent, we'll get a model that looks something like:



# Multi-Class Classification

Instead of there being two targets (pass/fail, cancer/not cancer, before/after 2000), we have  $K > 2$  targets.

Example:

- Beatles ( $K = 4$ ):
  - John Lennon, Paul McCartney, George Harrison, Ringo Starr
- Pets ( $K = \text{something large}$ ):
  - cat, dog, hamster, parrot, python, ...

# Representing the Targets

We use a **one-hot vector** to represent the target:

$$\mathbf{t} = (0, 0, \dots, 1, \dots, 0)$$

This vector contains  $K - 1$  zeros, and a single 1 somewhere.

Each *index* (column) in the vector represents one of the classes.

# Representing the Predictor

The prediction  $\mathbf{y}$  will also be a vector. Like in logistic regression there will be a linear part, and an activation function.

Linear part:  $\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$

So far, this is like having  $K$  separate logistic regression models, one for each element of the one-hot vector.

Q: What are the shapes of  $\mathbf{z}$ ,  $\mathbf{W}$ ,  $\mathbf{x}$  and  $\mathbf{b}$ ?

# Activation Function

Instead of using a *sigmoid* function, we instead use a **softmax activation** function:

$$y_k = \text{softmax}(z_1, \dots, z_K)_k = \frac{e^{z_k}}{\sum_{m=1}^K e^{z_m}}$$

The vector of predictions  $y_k$  is now a **probability distribution** over the classes!

# Why Softmax?

- Softmax is like the multi-class equivalent of sigmoid
- Softmax is a continuous analog of the “argmax” function
- If one of the  $z_k$  is much larger than the other, then the softmax will be approximately the argmax, in the one-hot encoding

# Cross-Entropy Loss

The cross-entropy loss naturally generalizes to the multi-class case:

$$\begin{aligned}\mathcal{L}(\mathbf{y}, \mathbf{t}) &= - \sum_{k=1}^K t_k \log(y_k) \\ &= -\mathbf{t}^\top \log(\mathbf{y})\end{aligned}$$

Recall that only one of the  $t_k$  is going to be 1, and the rest are 0.



# Multi-class Classification Summary

---

Model  $\mathbf{y} = \text{softmax}(\mathbf{W}^T \mathbf{x} + \mathbf{b})$

Loss  $\mathcal{L}(\mathbf{y}, \mathbf{t}) = -\mathbf{t}^T \log(\mathbf{y})$

Func-  
tion

Optimization  $\min_{\mathbf{w}, b} \mathcal{E}(\mathbf{w}, b)$  via Gradient Descent

Method

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \frac{\partial \mathcal{E}}{\partial \mathbf{W}}, \mathbf{b} \leftarrow \mathbf{b} - \alpha \frac{\partial \mathcal{E}}{\partial \mathbf{b}}$$

---

# Example: Beatle Recognition

Given a 100x100 pixel colour image of a face of a Beatle, identify the Beatle



Four possible labels:

- John Lennon
- Paul McCartney
- George Harrison
- Ringo Starr

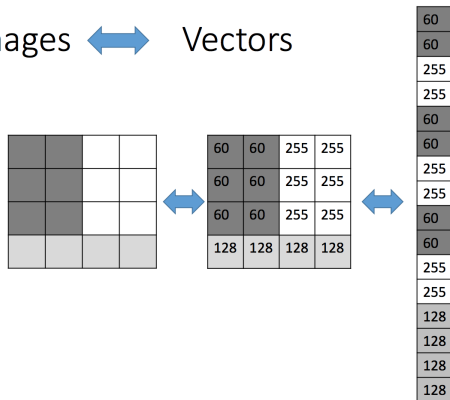
# Side Note: Representing an Image

This is what John Lennon looks like to a computer:

```
26 128 127 124 123 125 126 141 215 217 137 85 69 93 34 49 28 19 32 90 28 29 40 39 31 24 33 33 43 36 63 58 71 84 68 77
132 113 151 172 184 194 193 175 150 147 142 90 100 101 100 98 98 103 107 104 163 195 130 90 79 61 46 29 25 17 27 3
76 91 93 95 83 79 71 61 66 59 58 61 91 108 78 174 164 156 164 181 190 202 194 163 155 151 149 33 38 41 44 46 46 45
29 33 59 54 44 41 65 65 64 72 77 68 80 83 72 87 93 101 106 95 89 83 72 71 68 43 51 63 92 47 165 189 174 174 172 188 20
28 26 23 35 96 126 98 98 72 70 63 57 50 35 21 22 65 61 76 109 102 102 105 93 110 92 87 89 90 97 110 116 104 96 98 107
80 206 195 172 156 156 148 145 27 28 28 28 31 32 32 25 26 33 44 80 108 77 62 39 28 48 34 51 69 44 78 94 98 87 66 54
75 54 42 38 56 69 74 127 175 182 188 182 194 183 194 202 182 165 160 153 146 142 33 40 47 52 58 63 65 67 65 82 87 78 8
63 73 105 122 102 122 136 145 47 97 52 35 23 11 13 36 62 66 60 84 109 126 132 178 167 186 180 187 190 196 188 164 158
26 153 172 141 80 82 70 75 74 59 61 42 24 16 18 52 21 12 67 88 106 123 128 153 121 118 114 180 127 70 63 29 11 36 32 1
90 189 170 150 148 189 158 153 150 135 144 148 149 151 151 151 150 145 158 203 179 85 96 92 58 67 57 61 56 58 37 14 55
48 20 61 72 53 54 42 29 57 95 99 100 96 122 119 154 174 178 189 174 159 152 144 149 155 154 153 122 125 123 118 119 118
45 47 48 52 65 64 82 49 75 67 48 74 88 110 112 111 127 139 117 73 5 53 84 68 60 61 50 52 91 52 80 92 118 127 141 171
0 123 125 123 121 120 123 155 199 159 104 118 57 48 40 48 47 63 53 51 54 66 66 88 75 82 97 103 108 106 105 128 110 118
1 131 109 143 170 175 148 136 152 143 154 151 161 170 181 122 120 121 122 122 122 122 119 184 202 137 138 127 106
110 118 92 94 92 132 123 37 84 86 50 40 58 53 81 99 95 107 75 145 112 149 159 177 163 131 143 145 174 156 165 157 172
07 156 92 80 68 60 42 43 57 51 58 72 60 66 85 80 60 51 47 64 59 89 116 85 124 125 135 100 12 65 73 43 58 64 51 79 94 1
8 187 182 182 185 180 118 120 119 119 121 120 118 147 209 174 148 94 90 64 75 73 66 62 81 96 80 58 45 77 89 70 72 8
80 75 74 120 136 162 148 180 181 153 185 147 149 156 159 177 184 189 192 202 192 178 118 118 117 117 119 120 118 115 1
7 77 72 73 71 77 87 114 123 110 109 85 135 92 35 20 86 97 47 41 62 72 74 90 98 152 152 132 134 149 107 149 165 141 168
9 118 117 117 176 136 106 79 75 78 82 87 71 73 98 80 129 106 97 41 31 77 66 89 80 97 113 147 163 132 40 19 11 61 96 11
3 115 169 183 188 193 210 202 203 201 197 197 195 113 113 114 116 116 114 118 124 152 115 78 89 101 86 114 84 95 106 8
147 69 92 95 94 116 106 61 43 75 45 100 101 155 97 113 124 102 121 90 150 139 174 675 198 202 202 194 188 187 160 173
113 103 102 104 115 97 105 87 110 101 123 103 76 90 92 89 81 97 99 113 154 139 70 61 89 111 108 89 46 82 82 123 108 148
4 174 159 150 159 138 139 155 99 106 100 100 111 99 101 154 100 104 110 107 115 115 107 123 112 119 109 100 99 121 121
08 104 82 64 114 121 108 143 174 157 153 142 120 111 111 107 123 158 179 177 137 138 129 142 134 145 131 141 51 41
2 123 113 114 82 96 109 115 115 102 70 105 86 93 100 118 133 119 136 84 34 69 115 110 100 128 125 82 116 110 137 170 2
1 143 119 137 144 130 135 24 18 25 25 25 58 131 107 119 129 128 135 139 132 135 132 120 125 87 93 100 98 128 139 92
87 99 95 79 140 150 180 182 167 150 145 104 142 108 155 154 149 147 133 125 147 121 122 140 128 129 24 27 30 21 24 162
15 101 94 108 95 141 167 124 130 121 108 135 116 118 141 140 142 116 94 94 102 121 122 132 110 131 151 166 163 171 158
1 131 126 131 133 28 18 29 56 125 120 118 122 134 138 135 136 142 151 152 141 146 118 117 97 102 110 159 159 133 144
2 175 177 174 164 151 130 154 151 172 180 126 90 154 79 90 147 169 155 131 144 118 127 121 131 125 129 29 96 46 52 154
49 124 129 125 118 147 170 109 146 122 155 130 169 153 113 86 93 96 110 135 152 168 182 178 156 141 137 126 121 117 11
128 126 121 119 122 89 94 91 115 122 116 139 132 133 138 140 146 138 135 130 153 165 149 115 137 149 158 179 172 183
1 155 145 144 127 125 119 122 124 112 119 135 214 163 49 79 60 94 128 144 170 168 120 153 121 124 119 123 124 155 141
80 161 178 170 132 154 126 150 170 175 128 71 53 48 58 93 162 145 134 106 114 109 109 120 114 136 125 118 141 116 133 1
2 138 125 129 121 125 121 122 112 105 106 144 113 150 129 131 138 119 118 113 161 141 178 181 178 172 154 182 92 52 48
14 118 119 120 122 125 128 127 124 127 136 197 219 157 54 73 81 79 129 141 139 193 166 118 137 120 121 119 120 98 99 1
153 184 156 219 212 126 41 23 32 36 33 25 30 39 33 58 92 106 110 116 114 111 119 121 120 122 124 127 127 124 127 1
2 130 127 136 119 118 97 101 92 138 111 126 110 116 124 97 92 121 124 171 193 160 174 155 222 236 166 68 38 37 32 32
1 122 123 127 126 127 135 199 220 171 43 94 122 87 123 127 146 129 173 169 115 127 127 137 128 101 102 103 130 113
```

# Image as a Vector of Features

Images  $\longleftrightarrow$  Vectors



# Features and Targets

Each of our input images are  $100 \times 100$  pixels

$$\mathbf{y} = \text{softmax}(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

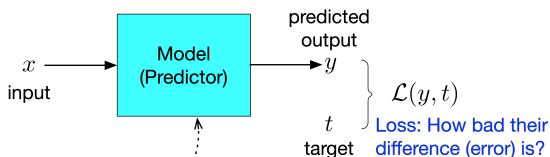
Questions:

- What will be the length of our input (feature) vectors  $\mathbf{x}$ ?
- What will be the length of our one-hot targets  $\mathbf{t}$ ?
- What are the shapes of  $\mathbf{W}$  and  $\mathbf{b}$ ?
- How many (scalar) parameters are in our model, in total?

## Section 4

# Deep Learning

# From Linear Models to Neural Networks



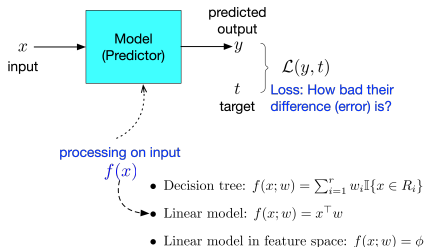
processing on input

$f(x)$

- Decision tree:  $f(x; w) = \sum_{i=1}^r w_i \mathbb{I}\{x \in R_i\}$
- Linear model:  $f(x; w) = x^\top w$
- Linear model in feature space:  $f(x; w) = \phi(x)^\top w$

- Design of an ML system follows a modular approach. We have to make choice on
  - Model
  - Loss function
  - Regularizer, etc.
- In your Introduction to ML course (ex. CSC311), you have seen many examples.

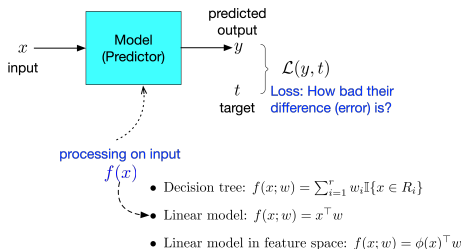
# Beyond Linear Models



- Feature mapping can make linear models much more powerful.
- Coming up with feature mapping can be challenging.
- Kernel-based approach is a way to partially address it.
- (Artificial) Neural Networks (NN) is a general approach to represent complex models.



# What is a Model After All?



- The predictor can be seen as a **computer program** that processes the input in order to generate the output. Some programs are simpler, some are more complex.
- Neural networks are general and flexible ways to specify a computer program.
- Different NN architectures correspond to different ways of specifying the overall architecture of the program.
- We are going to study several different ways in this course.

# What are Neural Networks? A Sneak Peak!

Neural networks are a class of models originally inspired by the brain.

- Most of the biological details aren't essential, so we use vastly simplified models of neurons.
- Nowadays we mostly think about math, statistics, etc
- Neural networks are collections of thousands (or millions) of these simple processing units that together perform useful computations

$$y = \phi(\mathbf{w}^T \mathbf{x} + \mathbf{b})$$

# Why Neural Networks?

- Very effective across a range of applications (vision, text, speech, medicine, robotics, etc.)
- Widely used in both academia and the tech industry
- Powerful software frameworks (PyTorch, TensorFlow, JAX, etc.) let us quickly implement sophisticated algorithms

# What is Deep Learning?

- A “deep” neural network contains many “layers”.
- Later layers use the output of earlier layers as input.
- The term **deep learning** emphasizes that the neural network algorithms often involve hierarchies with many stages of processing.

# Deep Learning Caveats: Interpretability

Before getting deep in studying NN and Deep Learning, it is good to know some of issues common with them.



Figure 5: from <https://xkcd.com/1838/>

# Deep Learning Caveats: Adversarial Examples



$x$

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

=



$x +$

$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

Image Credit: Goodfellow, Shlens, Szegedy, “Explaining and Harnessing Adversarial Examples,” ICLR, 2015.

**The U.S. military built an AI tool to find suitable combat personnel but had to shut it down because it was discriminating against women**

News



# Fairness in Machine Learning

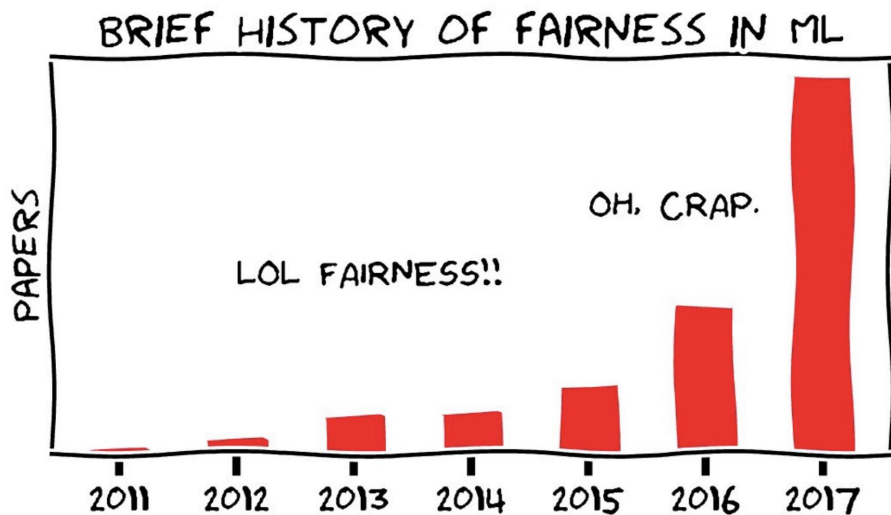


Image Credit: Solon Barocas and Moritz Hardt, "Fairness in Machine Learning", NeurIPS 2017



# Course Coverage

This is the tentative schedule and may slightly change.

- Mostly supervised learning
  - Multi-layer Perceptron (fully connected feedforward)
  - Convolutional Neural Networks for images
  - Recurrent Neural Networks for sequences
- Some unsupervised learning
  - Autoencoders
  - Generative Adversarial Networks
- Special topic TBD/TBA

## Section 5

# Logistics

# This Course

- This is the second course in machine learning, with a focus on neural networks and deep learning.
- First 75%: Mostly supervised learning
- Last 25%: Mostly unsupervised learning
- Three sections: equivalent content, different instructors, same deliverables.
- Only attend the section you are officially enrolled.
- Coursework is aimed at advanced undergrads. We will use multivariate calculus, probability, and linear algebra.

# What happens every week?

- We have three sessions per week.
- Each session (usually) consists of:
  - 2h of lecture
  - 1h of tutorial/practice (bring your laptop).
  - You start working on a homework/practice assignment during the tutorial/practice session. You have until next week to submit it.
  - We sometimes may use the tutorial time for lecture.

# How we communicate with you

- Course Website: **<https://amfarahmand.github.io/NN-Winter2024/>**
- Main source of information is the course webpage. Check regularly!
- We will also use Quercus for **announcements**.
- We will use Piazza for **discussions**.
- We will use MarkUs for **assignment submission**.

# Do I have the right background?

- **Machine Learning:** linear models for regression and classification, maximum likelihood estimation, PCA, EM, etc.
- **Linear algebra:** vector/matrix manipulations, basic properties of matrices.
- **Calculus:** partial derivatives/gradient.
- **Probability:** common distributions; Bayes' Rule.
- **Statistics:** expectation, variance, covariance, median; maximum likelihood.

# In the Classroom

- Feel free to ask questions from the instructors. Don't be shy! Asking questions helps you learn better. We try to answer as many questions as we can.
- Cell phones and other electronics are allowed in lecture (it might help your learning to annotate slides as we go through them).
- Talking with others is discouraged in the middle of lecture, as it would be distracting to us and others. We will have breaks so you can discuss among yourselves.
- Recording or taking pictures in class is **strictly prohibited** without the consent of your instructor. Please ask before doing!
- Even though pandemic is not in news anymore, COVID-19 and other respiratory diseases are still around. Some of your friends and peers might have weak immune system. Long COVID can be serious.
  - We encourage you to wear mask in the classroom and during in-person office hours

# Accomodations

- Please refer to <http://www.illnessverification.utoronto.ca> in case of illness (you need to fill out an absence declaration form on ACORN and contact me).
- If you require additional academic accommodations, please contact UofT Accessibility Services:  
<https://studentlife.utoronto.ca/department/accessibility-services/>



# Course Evaluation

This is tentative and may change in the next few days:

- Ten (10) practice assignments (40%)
  - Most (8) are small programming exercises.
  - Two (2) are derivation-based exercises.
  - You start working on them during the Tutorial part of each session. Our TAs will help you understand the assignment and guide you in solving them.
- Research Project (30%)
  - Research proposal (10%), written report and submitted codebase (20%).
- Take-Home Test: 20% (close to the end of the semester)
- Readings (10%)
  - Read some (5) research papers from a paper bank of 10-15 papers.
  - Write a short 1-paragraph summary and two questions on how the method(s) can be used or extended.

# Collaboration and Assignments

- Collaboration:
  - Collaboration on the Homework Assignments **is allowed**, under certain conditions:
    - You can discuss the assignment with another student (group of two).
    - In your submission, you need to be very clear about the contribution of each individual. For example, you should say we did a pair-programming or person A solved this part while person B solved another part.
    - You *can* use copilot, ChatGPT, etc. to solve the problem, but that would consider as your group member. That is, you can have either a 2-human group or 1 human + 1 machine group (no 2 machine group). If you do it, you should report it as well.
  - You need to form a team of 3–4 members to work on your projects (the exact number will be determined after finalizing the number of students enrolled).
    - Similar to the homework assignments, you need to report the contribution of each collaborator.
    - If you get the help of a machine, you need to clearly indicate that. The machine will cost you one of the team members.
  - Collaboration on the Take-home Test or Paper Readings is **not allowed**. These should be done as individual.

# Collaboration and Assignments

- Late Submissions (assignments, proposals, reports, etc)
  - Submissions should be handed in by deadline; a late penalty of 10% per day will be assessed thereafter (up to 3 days, then submission is blocked).
  - Extensions will be granted only in special situations, and you will need a Student Medical Certificate or a written request approved by the course coordinator at least one week before the due date.

# Academic Integrity

- By this point in your studies, you should know how to follow the academic integrity rules. You need to know what cheating and plagiarism are.
- If you need a review, read the U of T's Code of Behaviour on Academic Matters.
- Don't cheat or plagiarize!

# Useful Resources

Recommended readings will be given for each lecture. But the following will be useful throughout the course:

- **Deep Learning**, a textbook by Yoshua Bengio, Ian Goodfellow, and Aaron Courville.
- **Dive into Deep Learning**
- Video lectures for the U of T Professor Geoffrey Hinton's **course**.
- Andrej Karpathy's **lecture notes** on convolutional networks.
- Richard Socher's **lecture notes**, focusing on RNNs.
- Video lectures for Hugo Larochelle's **neural networks course**.
- If you need to brush up your basic knowledge of ML, you can take a look at one of the previous offerings of it at the U of T. For example, **CSC2515 - Fall 2022** (the content is almost the same as CSC311).

## Other Useful Resources:

- Trevor Hastie, Robert Tibshirani, and Jerome Friedman, The Elements of Statistical Learning, Second Edition, 2009.
- Christopher M. Bishop, Pattern Recognition and Machine Learning, 2006
- Richard S. Sutton and Andrew G. Barto, Reinforcement Learning: An Introduction, Second Edition, 2018.
- Amir-massoud Farahmand, Lecture Notes on Reinforcement Learning, 2021.
- Kevin Murphy, Machine Learning: A Probabilistic Perspective, 2012.
- Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani, An Introduction to Statistical Learning, 2017.
- Shai Shalev-Shwartz and Shai Ben-David, Understanding Machine Learning: From Theory to Algorithms, 2014.
- David MacKay, Information Theory, Inference, and Learning Algorithms, 2003.

# Compute

- Colaboratory: Programming assignments are to be completed in Google Colab, which is a web-based iPython Notebook service that has access to a free Nvidia K80 GPU per Google account. Highly recommended for homeworks and some course projects.
- Department Teaching Labs: Linux compute servers with desktop or datacentre-class GPUs. Recommended for course project.
- Google Compute Engine: GCE delivers virtual machines running in Google's data center. Recommended for course project.

## Section 6

What to do this week?



# What to do this week?

- No tutorials this week.
- Review your linear algebra, probability, and ML background