





## Advanced Multimodal Machine Learning

**Lecture 2.1: Basic Concepts** 

**Louis-Philippe Morency Tadas Baltrusaitis** 

#### **Lecture Objectives**

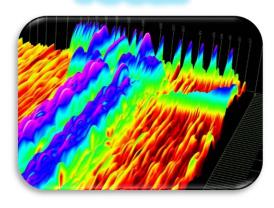
- Unimodal basic representations
  - Visual, language and acoustic modalities
- Data-driven machine learning
  - Training, validation and testing
  - Example: K-nearest neighbor
- Linear Classification
  - Score function
  - Two loss functions (cross-entropy and hinge loss)
- Course project "speed-dating"

#### **Multimodal Machine Learning**

## Verbal



## Vocal



## Visual



## **Core Technical Challenges:**

Representation Translation

**Alignment** 

Fusion

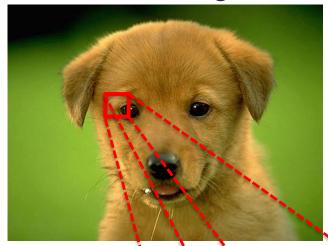
Co-Learning

These challenges are non-exclusive.



# Unimodal Basic Representations

#### Color image



Each pixel is represented in  $\mathcal{R}^d$ , d is the number of colors (d=3 for RGB)

1									
	88	82	84	88	85	83	80	93	102
	88	80	78	80	80	78	73	94	100
	85	79	80	78	77	74	65	91	99
	38	35	40	35	39	74	77	70	65
	20	25	23	28	37	69	64	60	57
	22	26	22	28	40	65	64	59	34
	24	28	24	30	37	60	58	56	66
į	21	22	23	27	38	60	67	65	67
i	23	22	22	25	38	59	64	67	66

Input observation  $x_i$ 

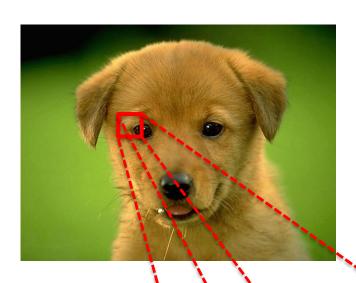
#### 

# Binary classification problem



Dog?

label  $y_i \in \mathcal{Y} = \{0,1\}$ 



Each pixel is represented in  $\mathcal{R}^d$ , d is the number of colors (d=3 for RGB)

1											
	88	82	84	88	85	83	80	93	102		
	88	80	78	80	80	78	73	94	100		
	85	79	80	78	77	74	65	91	99		
	38	35	40	35	39	74	77	70	65		
	20	25	23	28	37	69	64	60	57		
	22	26	22	28	40	65	64	59	34		
	24	28	24	30	37	60	58	56	66		
Ì	21	22	23	27	38	60	67	65	67		
Ì	23	22	22	25	38	59	64	67	66		

Input observation  $x_i$ 

21 23

82

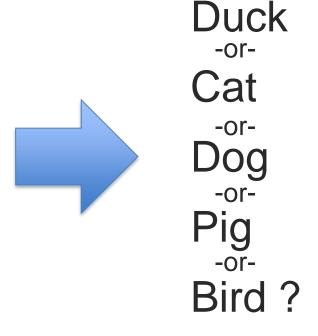
80

26

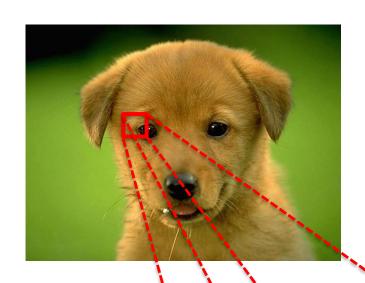
78

80

# Multi-class classification problem



label  $y_i \in \mathcal{Y} = \{0,1,2,3,...\}$ 



Each pixel is represented in  $\mathcal{R}^d$ , d is the number of colors (d=3 for RGB)

1		1							
	88	82	84	88	85	83	80	93	102
	88	80	78	80	80	78	73	94	100
	85	79	80	78	77	74	65	91	99
	38	35	40	35	39	74	77	70	65
	20	25	23	28	37	69	64	60	57
	22	26	22	28	40	65	64	59	34
	24	28	24	30	37	60	58	56	66
į	21	22	23	27	38	60	67	65	67
Ì	23	22	22	25	38	59	64	67	66

nput observation  $x_i$ 

21 23 82

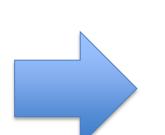
80

26

28

80

# Multi-label (or multi-task) classification problem



Duck?

Cat?

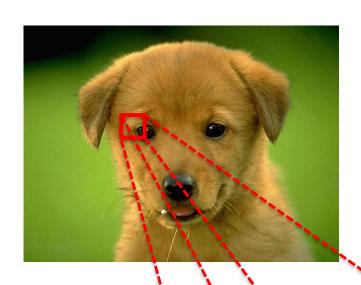
Dog?

Pig?

Bird?

Puppy?

label vector  $\mathbf{y_i} \in \mathcal{Y}^m = \{0,1\}^m$ 



Each pixel is represented in  $\mathbb{R}^d$ , d is the number of colors (d=3 for RGB)

88	82	84	88	85	83	80	93	102
88	80	78	80	80	78	73	94	100
85	79	80	78	77	74	65	91	99
38	35	40	35	39	74	77	70	65
20	25	23	28	37	69	64	60	57
22	26	22	28	40	65	64	59	34
24	28	24	30	37	60	58	56	66
21	22	23	27	38	60	67	65	67
23	22	22	25	38	59	64	67	66
	88 85 38 20 22 24 21	88 80 85 79 38 35 20 25 22 26 24 28 21 22	88 80 78 85 79 80 38 35 40 20 25 23 22 26 22 24 28 24 21 22 23	88 80 78 80 85 79 80 78 38 35 40 35 20 25 23 28 22 26 22 28 24 28 24 30 21 22 23 27	88     80     78     80     80       85     79     80     78     77       38     35     40     35     39       20     25     23     28     37       22     26     22     28     40       24     28     24     30     37       21     22     23     27     38	88     80     78     80     80     78       85     79     80     78     77     74       38     35     40     35     39     74       20     25     23     28     37     69       22     26     22     28     40     65       24     28     24     30     37     60       21     22     23     27     38     60	88     80     78     80     80     78     73       85     79     80     78     77     74     65       38     35     40     35     39     74     77       20     25     23     28     37     69     64       22     26     22     28     40     65     64       24     28     24     30     37     60     58       21     22     23     27     38     60     67	88     80     78     80     80     78     73     94       85     79     80     78     77     74     65     91       38     35     40     35     39     74     77     70       20     25     23     28     37     69     64     60       22     26     22     28     40     65     64     59       24     28     24     30     37     60     58     56       21     22     23     27     38     60     67     65

nput observation  $x_i$ 

21 23

82

80

35

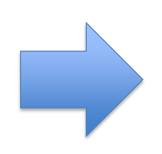
26

28

22

80

# Multi-label (or multi-task) regression problem



Age?

Height?

Weight?

Distance?

Happy?

label vector  $y_i \in \mathcal{Y}^m = \mathbb{R}^m$ 

#### **Unimodal Classification – Language Modality**



#### Masterful!

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful

#### MARTHA (CON'T)

Look around you. Look at all the great things you've done and the people you've helped.

#### CLARK

But you've only put up the good things they say about me.

#### MARTHA

Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

Input observation  $oldsymbol{x}$ 

## Word-level classification

Part-of-speech?



Named entity? (names of person,...)



"one-hot" vector

 $|x_i|$  = number of words in dictionary



#### **Unimodal Classification – Language Modality**



#### Masterful!

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful

#### MARTHA (CON'T)

Look around you. Look at all the great things you've done and the people you've helped.

#### CLARK

But you've only put up the good things they say about me.

#### MARTHA

Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

Input observation x

## Document-level classification



"bag-of-word" vector

 $|x_i|$  = number of words in dictionary



## **Unimodal Classification – Language Modality**



#### Masterful!

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful

#### MARTHA (CON'T)

Look around you. Look at all the great things you've done and the people you've helped.

#### OT ADV

But you've only put up the good things they say about me.

#### MARTHA

Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

Input observation  $x_i$ 

## Utterance-level classification

Sentiment? (positive or negative)

#### "bag-of-word" vector

 $|x_i|$  = number of words in dictionary

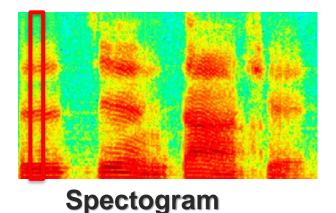


#### **Unimodal Classification – Acoustic Modality**

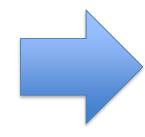
#### Digitalized acoustic signal



- Sampling rates: 8~96kHz
- Bit depth: 8, 16 or 24 bits
- Time window size: 20ms
  - Offset: 10ms



 $\begin{array}{c} \mathbf{n} \\ \mathbf{$ 



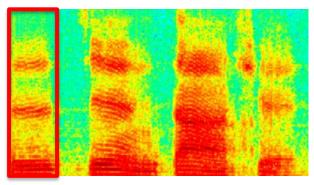
Spoken word?

#### **Unimodal Classification – Acoustic Modality**

#### Digitalized acoustic signal



- Sampling rates: 8~96kHz
- Bit depth: 8, 16 or 24 bits
- Time window size: 20ms
  - Offset: 10ms



**Spectogram** 





**Emotion?** 

Spoken word?

Voice quality?

# Data-Driven Machine Learning

#### **Data-Driven Machine Learning**

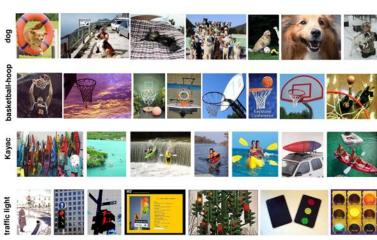
- **1. Dataset:** Collection of labeled samples D:  $\{x_i, y_i\}$
- 2. Training: Learn classifier on training set
- 3. Testing: Evaluate classifier on hold-out test set



#### Simple Classifier?







Dataset

## Traffic light

or-

-or-

**Basket** 

-or-

Kayak?

#### Simple Classifier: Nearest Neighbor







Traffic light

-or-

Dog

-or-

**Basket** 

-or-

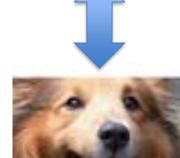
Kayak?

#### **Nearest Neighbor Classifier**

- Non-parametric approaches—key ideas:
  - "Let the data speak for themselves"
  - "Predict new cases based on similar cases"
  - "Use multiple local models instead of a single global model"
- What is the complexity of the NN classifier w.r.t training set of N images and test set of M images?
  - at training time?O(1)
  - At test time?
    O(N)

## Simple Classifier: Nearest Neighbor





#### **Distance metrics**

L1 (Manhattan) distance:

$$d_1(x_1, x_2) = \sum_{j} \left| x_1^{j} - x_2^{j} \right|$$

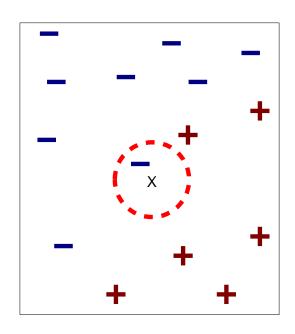
L2 (Eucledian) distance:

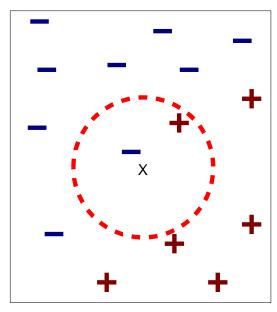
$$d_2(x_1, x_2) = \sqrt{\sum_j \left(x_1^j - x_2^j\right)^2}$$

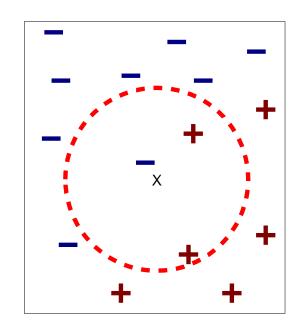
Which distance metric to use?

First hyper-parameter!

## **Definition of K-Nearest Neighbor**







- (a) 1-nearest neighbor
- (b) 2-nearest neighbor
- (c) 3-nearest neighbor

#### What value should we set K?

Second hyper-parameter!

#### **Data-Driven Approach**

- **1. Dataset:** Collection of labeled samples D:  $\{x_i, y_i\}$
- 2. Training: Learn classifier on training set
- 3. Validation: Select optimal hyper-parameters
- 4. Testing: Evaluate classifier on hold-out test set

Training Data

Validation Data

Test Data

#### **Evaluation methods (for validation and testing)**

- Holdout set: The available data set D is divided into two disjoint subsets,
  - the *training set D<sub>train</sub>* (for learning a model)
  - the  $test set D_{test}$  (for testing the model)
- Important: training set should not be used in testing and the test set should not be used in learning.
  - Unseen test set provides a unbiased estimate of accuracy.
- The test set is also called the holdout set. (the examples in the original data set D are all labeled with classes.)
- This method is mainly used when the data set D is large.
- Holdout methods can also be used for validation

#### **Evaluation methods (for validation and testing)**

- n-fold cross-validation: The available data is partitioned into n equal-size disjoint subsets.
- Use each subset as the test set and combine the rest n-1 subsets as the training set to learn a classifier.
- The procedure is run *n* times, which give *n* accuracies.
- The final estimated accuracy of learning is the average of the n accuracies.
- 10-fold and 5-fold cross-validations are commonly used.
- This method is used when the available data is not large.

#### **Evaluation methods (for validation and testing)**

- Leave-one-out cross-validation: This method is used when the data set is very small.
- Each fold of the cross validation has only a single test example and all the rest of the data is used in training.
- If the original data has m examples, this is mfold cross-validation
- It is a special case of cross-validation

# Linear Classification: Scores and Loss

#### Linear Classification (e.g., neural network)

#### **I**mage









- 1. Define a (linear) score function
- 2. Define the loss function (possibly nonlinear)
- 3. Optimization

#### 1) Score Function







Duck?

Cat?

Dog?

Pig?

Bird?

What should be the prediction score for each label class?

#### For linear classifier:

 $f(x_i; W, b) = Wx_i + b$ Weights [10x3072] Bias vector

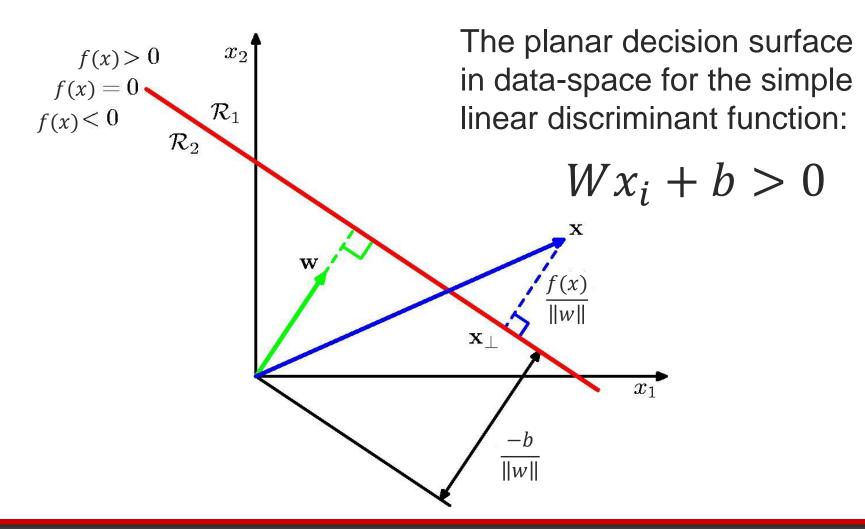
Parameters [10x3073]

**Input observation** (*i*<sup>th</sup> element of the dataset)

[3072x1]

[10x1]

#### **Interpreting a Linear Classifier**



#### Some Notation Tricks – Multi-Label Classification

$$W = \begin{bmatrix} W_1 & W_2 & \dots & W_N \end{bmatrix}$$

$$f(x_i; W, b) = Wx_i + b$$

$$f(x_i; W) = Wx_i$$

Weights x Input + Bias

[10x3072] [3072x1] [10x1]

Weights x Input

[10x3073] [3073x1]

The bias vector will be the last column of the weight matrix

Add a "1" at the end of the input observation vector

#### **Some Notation Tricks**

General formulation of linear classifier:

$$f(x_i; W, b)$$

"dog" linear classifier:

$$f(x_i; W_{dog}, b_{dog})$$
 or

$$f(x_i; W, b)_{dog}$$

or

 $f_{dog}$ 

Linear classifier for label *j*:

$$f(x_i; W_j, b_j)$$

or

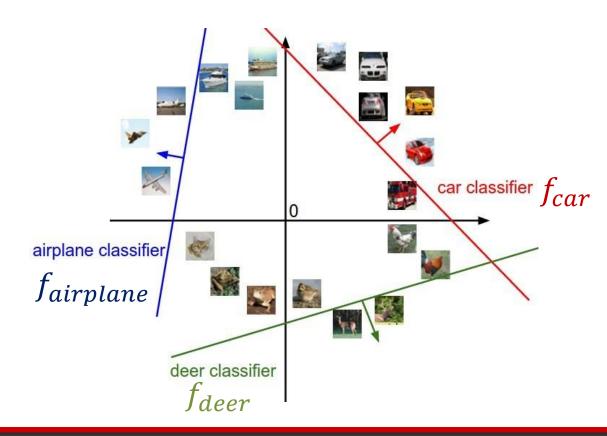
$$f(x_i; W, b)_j$$

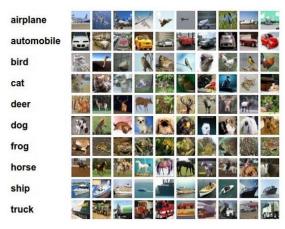
or

 $f_{j}$ 

#### **Interpreting Multiple Linear Classifiers**

$$f(x_i; W_j, b_j) = W_j x_i + b_j$$

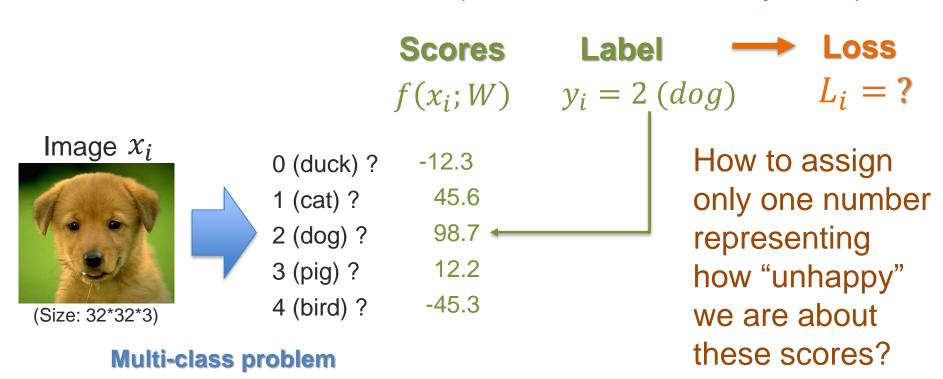




CIFAR-10 object recognition dataset

#### **Linear Classification: 2) Loss Function**

(or cost function or objective)



The loss function quantifies the amount by which the prediction scores deviate from the actual values.

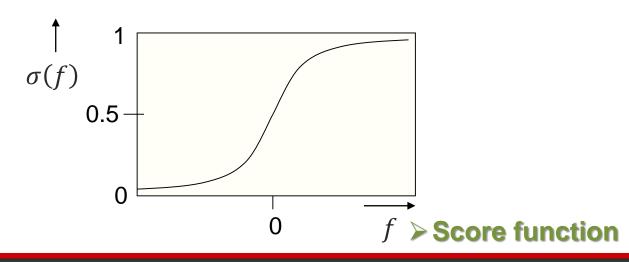


A first challenge: how to normalize the scores?

(or logistic loss)

Logistic function:

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



(or logistic loss)

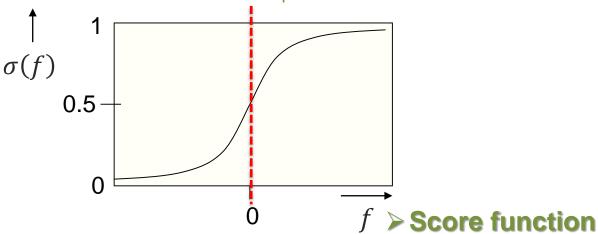
Logistic function:

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

Logistic regression: (two classes)

$$p(y_i = "dog"|x_i; w) = \sigma(w^T x_i)$$
= true

for two-class problem



(or logistic loss)

Logistic function:

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

Logistic regression: (two classes)

$$p(y_i = "dog"|x_i; w) = \sigma(w^T x_i)$$
= true

for two-class problem

Softmax function: (multiple classes)

$$p(y_i|x_i;W) = \frac{e^{f_{y_i}}}{\sum_j e^{f_j}}$$

Cross-entropy loss:

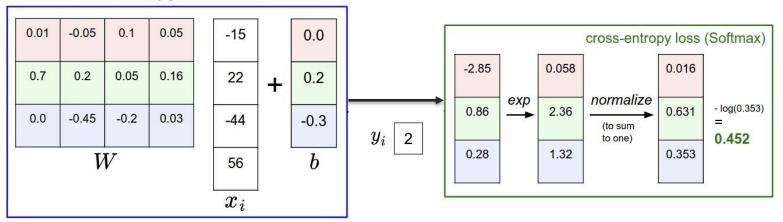
(or logistic loss)

$$L_i = -\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}}\right)$$

Softmax function

Minimizing the negative log likelihood.

matrix multiply + bias offset



#### **Second Loss Function: Hinge Loss**

(or max-margin loss or Multi-class SVM loss)

$$L_i = \sum_{j 
eq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$
 loss due to

example i sum over all incorrect labels

difference between the correct class score and incorrect class score



#### **Second Loss Function: Hinge Loss**

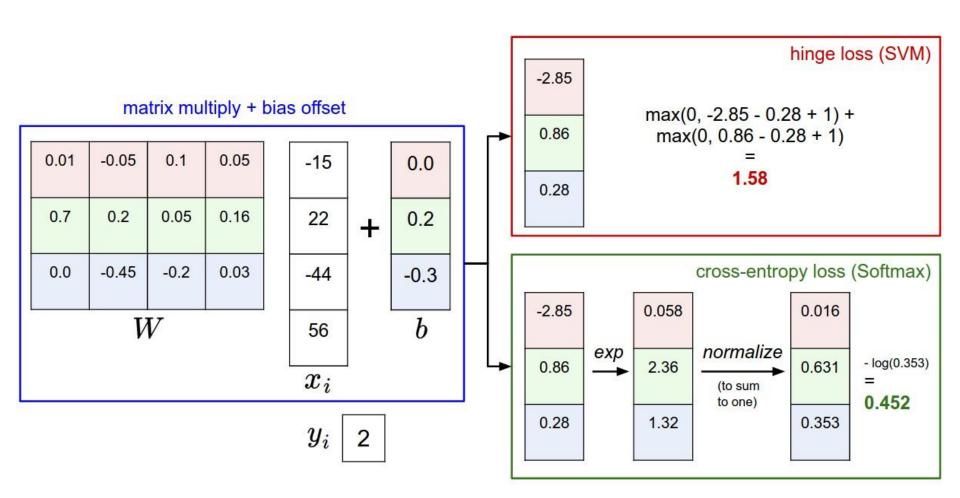
(or max-margin loss or Multi-class SVM loss)

$$L_i = \sum_{j 
eq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + extstyle{\Delta})$$
 e.g. 10

Example: 
$$f(x_i,W) = [13,-7,11] \ y_i = 0$$

$$L_i = \max(0, -7 - 13 + 10) + \max(0, 11 - 13 + 10)$$

#### **Two Loss Functions**





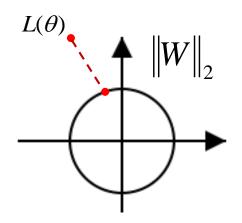
How to find the optimal W?

#### Regularization

$$L_{i} = -\log\left(\frac{e^{f_{y_{i}}(x_{i};W)}}{\sum_{j} e^{f_{j}(x_{i};W)}}\right) + \lambda R(W)$$
Regularization factor

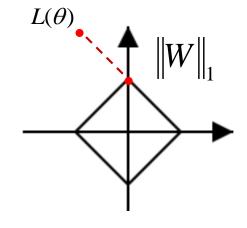
#### L-2 Norm (Gaussian prior):

$$R(W) = ||W||_2$$



#### L-1 Norm (Laplacian prior):

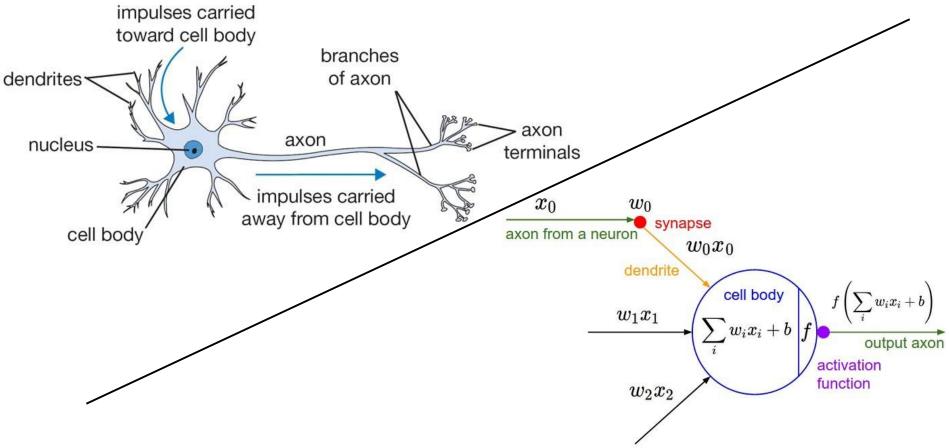
$$R(W) = ||W||_1$$



# Basic Concepts: Neural Networks

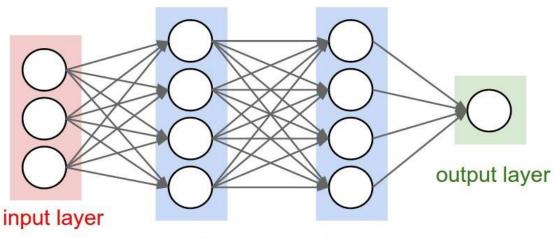
## **Neural Networks – inspiration**

Made up of artificial neurons



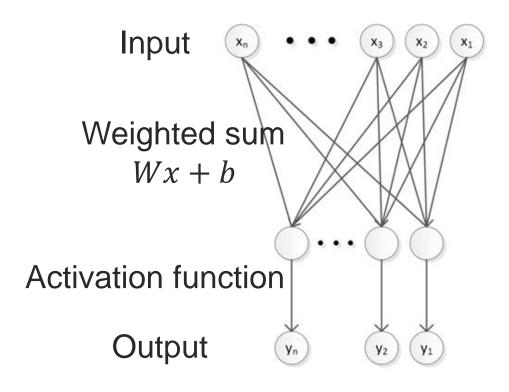
#### **Neural Networks – score function**

- Made up of artificial neurons
  - Linear function (dot product) followed by a nonlinear activation function
- Example a Multi Layer Perceptron



## **Basic NN building block**

Weighted sum followed by an activation function



$$y = f(Wx + b)$$

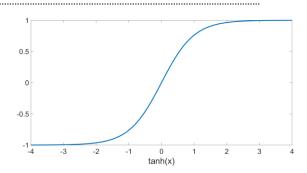
#### **Neural Networks – activation function**

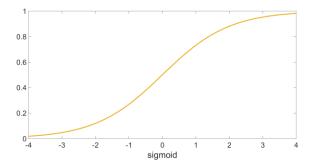
• 
$$f(x) = \tanh(x)$$

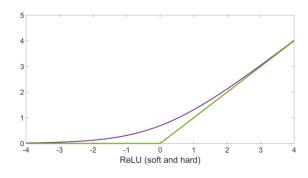
• Sigmoid - 
$$f(x) = (1 + e^{-x})^{-1}$$

• Linear 
$$-f(x) = ax + b$$

- ReLU  $f(x) = \max(0, x) \sim \log(1 + \exp(x))$ 
  - Rectifier Linear Units
  - Faster training no gradient vanishing
  - Induces sparsity







#### **Neural Networks – loss function**

- Already discussed it cross-entropy, Euclidean loss, cosine similarity, etc.
- Combine it with the score function to have an end-to-end training objective
- As example use Euclidean loss for data-point /  $L_i = (f(x_i) y_i)^2 = (f_{3;W_3}(f_{2;W_2}(f_{1;W_1}(x_i)))^2$
- Full loss is computed across all training samples

$$L = \sum_{i} (f(x_i) - y_i)^2$$

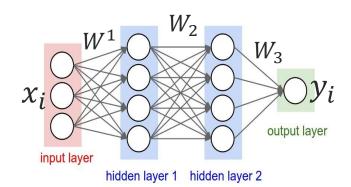
#### **Multi-Layer Feedforward Network**

#### Activation functions (individual layers)

$$f_{1;W_1}(x) = \sigma(W_1 x + b_1)$$

$$f_{2;W_2}(x) = \sigma(W_2x + b_2)$$

$$f_{3;W_3}(x) = \sigma(W_3 x + b_3)$$



#### Score function

$$y_i = f(x_i) = f_{3;W_3}(f_{2;W_2}(f_{1;W_1}(x_i)))$$

#### Loss function (e.g., Euclidean loss)

$$L_i = (f(x_i) - y_i)^2 = (f_{3;W_3}(f_{2;W_2}(f_{1;W_1}(x_i))))^2$$