Topics:

- Linear Classification, Loss functions
- Gradient Descent

# CS 4644-DL / 7643-A
# ZSOLT KIRA

- **Assignment 1 out!**
  - Due **Feb. 2$^{nd}$ 11:59pm** (grace period Feb 4$^{th}$).
  - Start early, start early, start early!
  - HW1 Tutorial, Matrix Calculus tutorial OH: TBA

- **Piazza**: Should be sync'd with Canvas now (Code: DLSPR2024)
  - **NOTE:** There is an OMSCS section with a Ed. Make sure you are in the right one

- **Tentative Office hours** schedule:
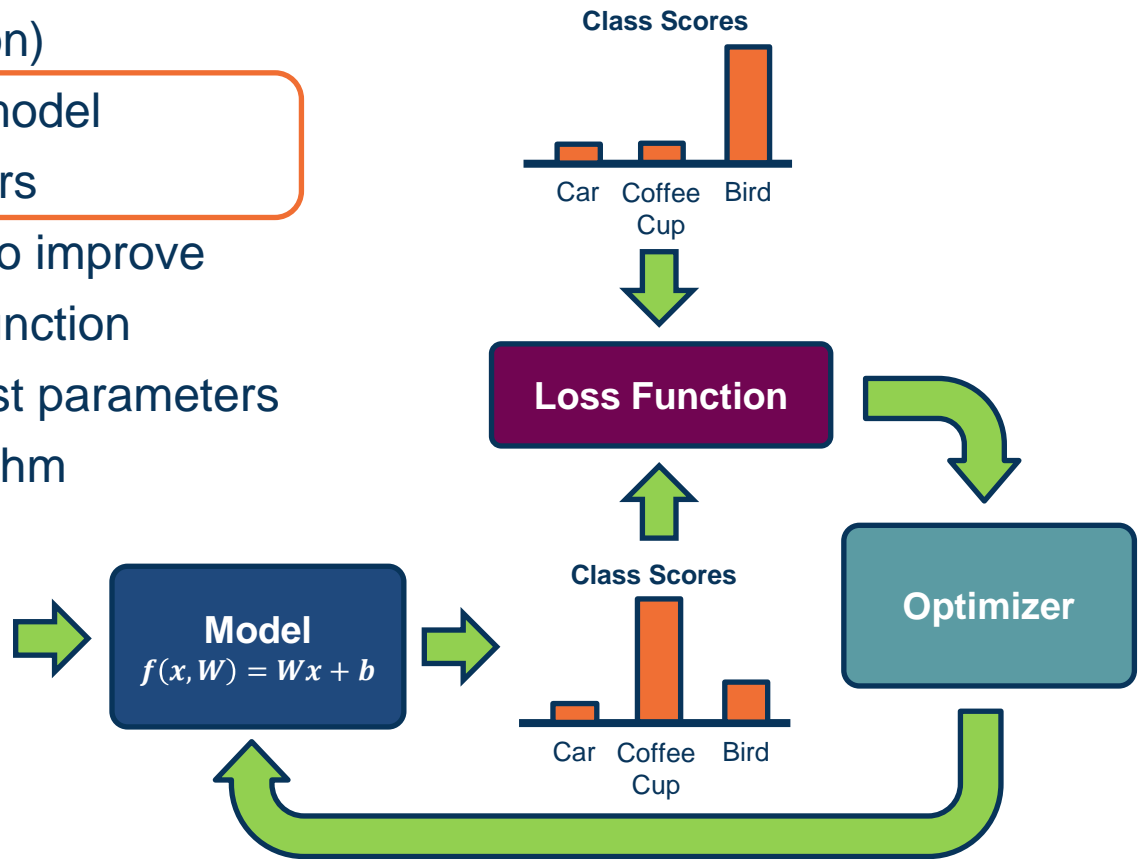https://piazza.com/class/lcl94yjxkbb59e/post/53
  - Calendar on webpage:
https://faculty.cc.gatech.edu/~zk15/teaching/AY2024_cs7643_spring/index.html
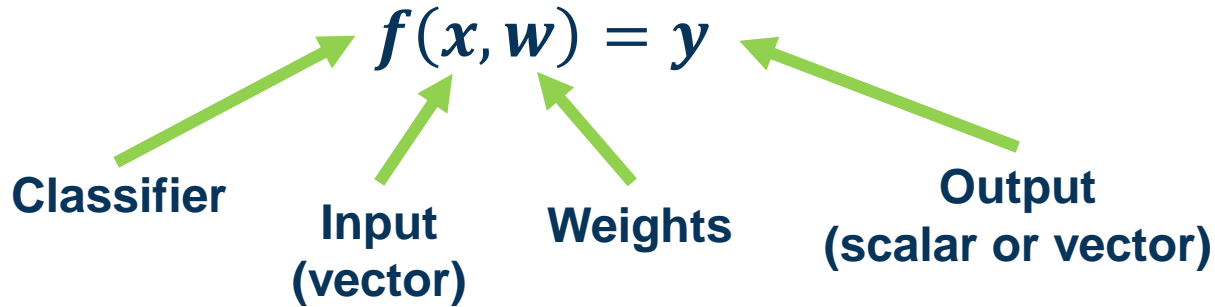
- Input (and representation)
- Functional form of the model
  - Including parameters
- Performance measure to improve
  - Loss or objective function
- Algorithm for finding best parameters
  - Optimization algorithm

Class Scores
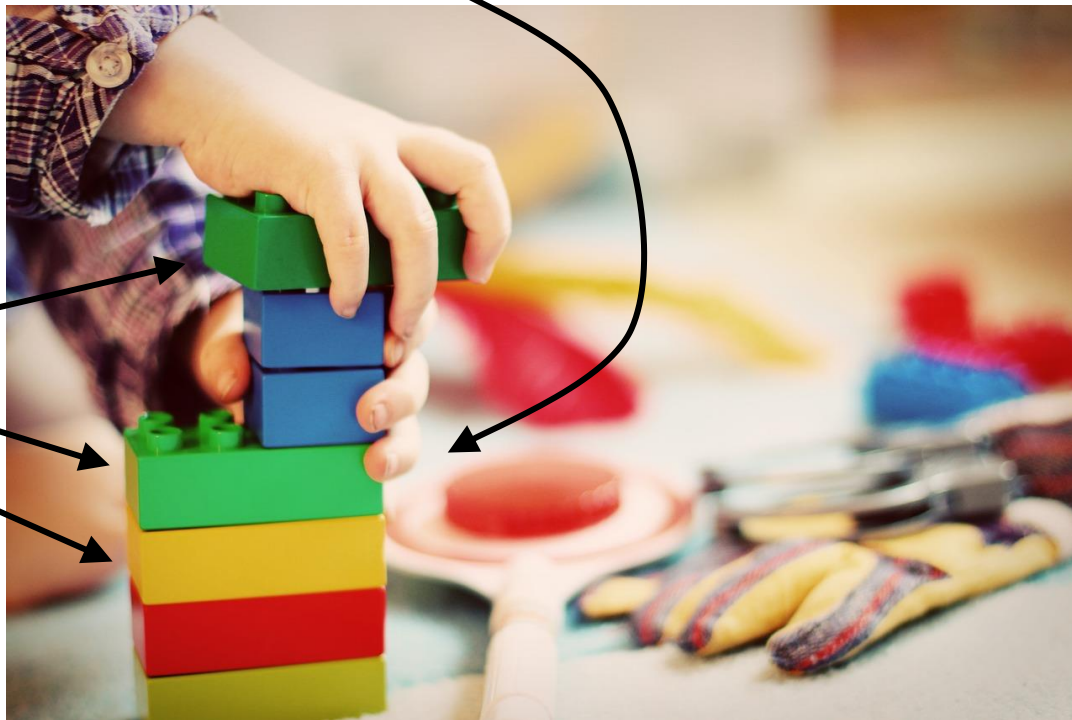
Car    Coffee    Bird
        Cup

**Loss Function**

**Optimizer**

**Model**
$f(x, W) = Wx + b$

**Data: Image**

Class Scores

Car    Coffee    Bird
        Cup

**Components of a Parametric Model**

Georgia
Tech

$$f(x, w) = y$$

**Classifier**

**Input (vector)**

**Weights**

**Output (scalar or vector)**

- **Input:** Continuous number or vector
- **Output:** A continuous number
  - For classification typically a **score**
  - For regression what we want to regress to (house prices, crime rate, etc.)
- $w$ **is a vector and weights** to optimize to fit target function

**Model: Discriminative Parameterized Function**

Georgia Tech

Neural Network

Linear classifiers

This image is CC0 1.0 public domain

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

**Deep Learning as Legos**

Georgia Tech

- We can move the bias term into the weight matrix, and a "1" at the end of the input

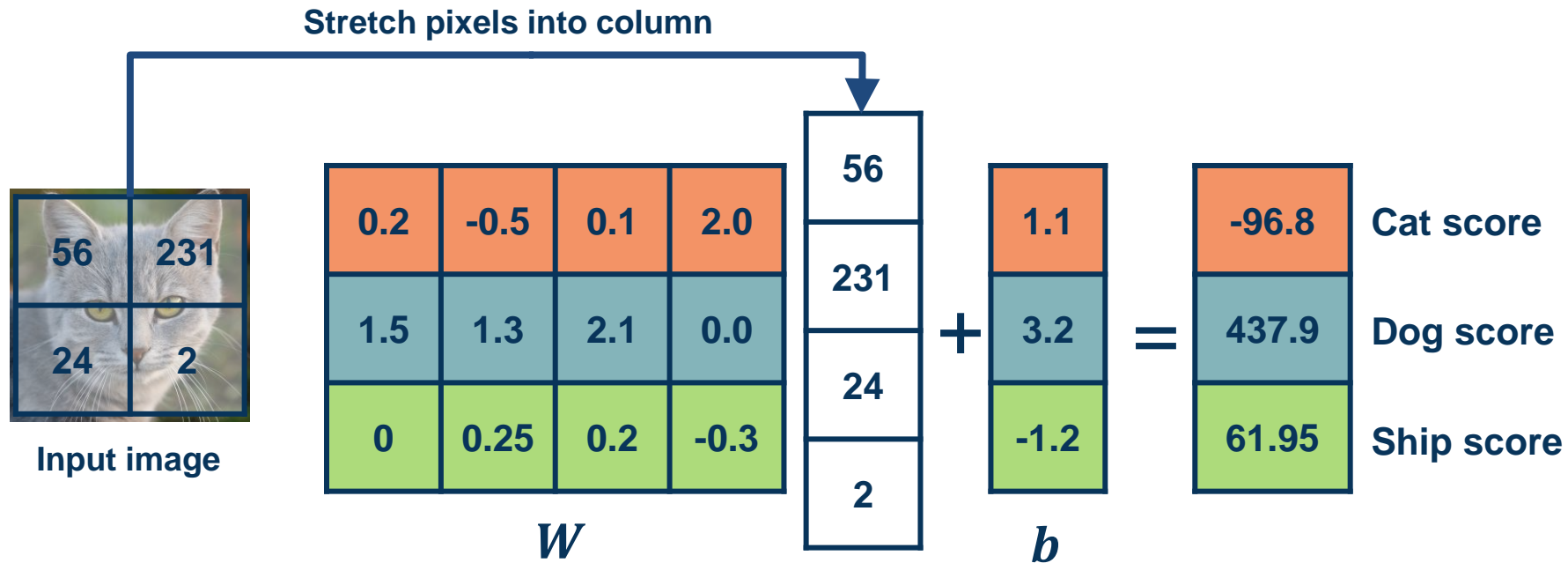- Results in **one matrix-vector multiplication**!

**Model**

$$f(x, W) = Wx + b$$

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b_1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b_2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \\ 1 \end{bmatrix}$$
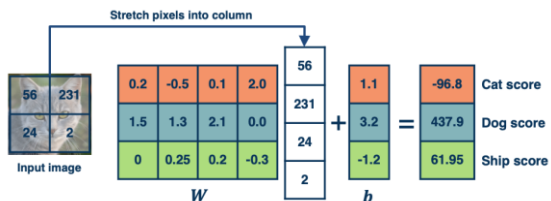
$$W \qquad\qquad x$$

Georgia Tech

# Example with an image with **4 pixels**, and **3 classes (cat/dog/ship)**

**Stretch pixels into column**

| | | | |
|---|---|---|---|
| 0.2 | -0.5 | 0.1 | 2.0 |
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

| |
|---|
| 56 |
| 231 |
| 24 |
| 2 |

$W$

**+**

| |
|---|
| 1.1 |
| 3.2 |
| -1.2 |

$b$

**=**

| | |
|---|---|
| -96.8 | **Cat score** |
| 437.9 | **Dog score** |
| 61.95 | **Ship score** |

**Input image**

56    231

24    2

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**Example**

Georgia Tech

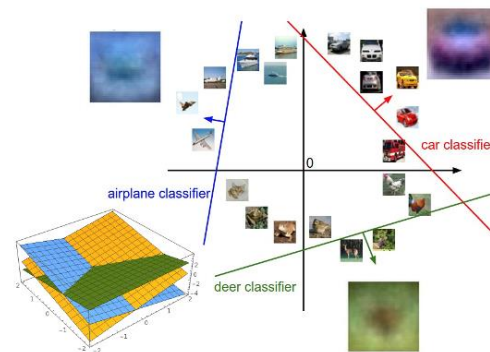| Algebraic Viewpoint | Visual Viewpoint | Geometric Viewpoint |
|---|---|---|
| $f(x, W) = Wx$ | One template per class | Hyperplanes cutting up space |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*
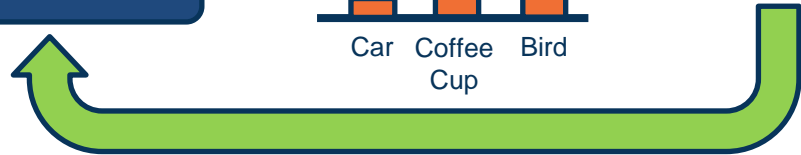
**Linear Classifier: Three Viewpoints**

Georgia Tech
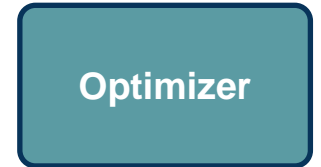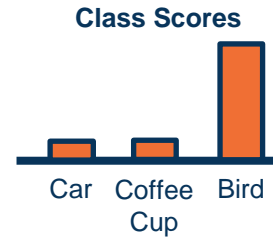
Performance Measure for a Classifier

- Input (and representation)
- Functional form of the model
  - Including parameters
- **Performance measure to improve**
  - **Loss or objective function**
- Algorithm for finding best parameters
  - Optimization algorithm

**Data: Image**

**Features: Histogram**

**Model**
$f(x, W) = W_x + b$

Class Scores

Car   Coffee   Bird
       Cup

Class Scores

Car   Coffee   Bird
       Cup

**Loss Function**

**Optimizer**

**Components of a Parametric Model**

Georgia Tech

The output of a classifier can be considered a **score**

For binary classifier, use rule:

$$y = \begin{cases} 1 & \textbf{if } f(x, w) >= 0 \\ 0 & \textbf{otherwise} \end{cases}$$

- Can be used for many classes by considering one class versus all the rest (one versus all)

For multi-class classifier can take the maximum

**Class Scores**

Model
$f(x, W) = Wx + b$

Car     Coffee     Bird
        Cup

Several issues with scores:

- Not very interpretable (no bounded value)

We often want **probabilities**

- More interpretable

- Can relate to probabilistic view of machine learning

We use the **softmax** function to convert scores to probabilities

$$s = f(x, W)$$ **Scores**

$$P(Y = k | X = x) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ **Softmax Function**

**Converting Scores to Probabilities**

Georgia Tech

We need a performance measure to **optimize**

- Penalizes model for being wrong

- Allows us to modify the model to reduce this penalty

- Known as an **objective** or **loss** function

In machine learning we use **empirical risk minimization**

- Reduce the loss over the **training** dataset

- We **average** the loss over the training data

Given a dataset of examples**:**

$$\{(x_i, y_i)\}_{i=1}^N$$

Where $x_i$ is image and

$y_i$ is (integer) label

Loss over the dataset is a sum of loss over examples:

$$L = \frac{1}{N} \sum L(f(x_i, W), y_i)$$

# Multiclass SVM loss:

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

**margin**

**scores for other classes**  **score for correct class**

**score**

**Example: "Hinge Loss"**

$s_{y_i}$

$s_j$  1

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**Performance Measure for Scores**

Georgia Tech

## Multiclass SVM loss:

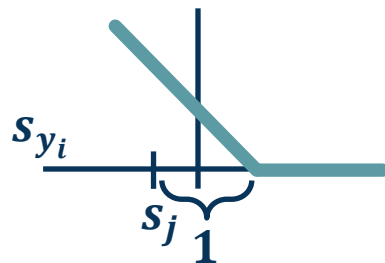Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

= max(0, 5.1 - 3.2 + 1)
+max(0, -1.7 - 3.2 + 1)
= max(0, 2.9) + max(0, -3.9)
= 2.9 + 0
= 2.9

Suppose: 3 training examples, 3 classes. With some $W$ the scores $f(x,W)=Wx$ are:



| | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| **Losses:** | **2.9** | | |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**SVM Loss Example**

Georgia Tech

## Multiclass SVM loss:

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

= max(0, 1.3 - 4.9 + 1)
    +max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0
= 0

Suppose: 3 training examples, 3 classes.
With some $W$ the scores $f(x,W)=Wx$ are:



|       | cat  | car  | frog |
|-------|------|------|------|
| cat   | 3.2  | 1.3  | 2.2  |
| car   | 5.1  | 4.9  | 2.5  |
| frog  | -1.7 | 2.0  | -3.1 |
| Losses: |    | 0.0  |      |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

## SVM Loss Example

Georgia Tech

**Multiclass SVM loss:**

Suppose: 3 training examples, 3 classes.
With some $W$ the scores $f(x, W) = Wx$ are:

$$L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

Q: What happens to loss if car image scores change a bit?

No change for small values

|       |       |       |       |
|-------|-------|-------|-------|
| cat   | **3.2** | 1.3   | 2.2   |
| car   | 5.1   | **4.9** | 2.5   |
| frog  | -1.7  | 2.0   | **-3.1** |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**SVM Loss Example**

Georgia Tech

**Multiclass SVM loss:**

$$L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

Q: What is min/max of loss value?

[0,inf]

Suppose: 3 training examples, 3 classes.
With some $W$ the scores $f(x, W) = Wx$ are:



|       |        |        |         |
|-------|--------|--------|---------|
| cat   | **3.2**| 1.3    | 2.2     |
| car   | 5.1    | **4.9**| 2.5     |
| frog  | -1.7   | 2.0    | **-3.1**|

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**SVM Loss Example**

Georgia Tech

**Multiclass SVM loss:**

$$L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

Q: At initialization W is small so all s ≈ 0.
What is the loss?

C-1

Suppose: 3 training examples, 3 classes.
With some $W$ the scores $f(x,W)=Wx$ are:



|      | cat  | car  | frog |
|------|------|------|------|
| cat  | **3.2**  | 1.3  | 2.2  |
| car  | 5.1  | **4.9**  | 2.5  |
| frog | -1.7 | 2.0  | **-3.1** |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**SVM Loss Example**

Georgia Tech

**Multiclass SVM loss:**

Suppose: 3 training examples, 3 classes.
With some $W$ the scores $f(x,W)=Wx$ are:

$$L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

Q: What if the sum was over all classes? (including j = y_i)

No difference (add constant 1)



|      |      |      |      |
|------|------|------|------|
| cat  | **3.2** | 1.3  | 2.2  |
| car  | 5.1  | **4.9** | 2.5  |
| frog | -1.7 | 2.0  | **-3.1** |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**SVM Loss Example**

Georgia Tech

**Multiclass SVM loss:**

$$L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

Q: What if we used mean instead of sum?

No difference
Scaling by constant

Suppose: 3 training examples, 3 classes.
With some $W$ the scores $f(x,W)=Wx$ are:



| | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**SVM Loss Example**

Georgia Tech

## Multiclass SVM loss:

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$$

L = (2.9 + 0 + 12.9)/3
   = **5.27**

Suppose: 3 training examples, 3 classes.
With some $W$ the scores $f(x,W)=Wx$ are:



|       | cat | car | frog |
|-------|-----|-----|------|
| cat   | **3.2** | 1.3 | 2.2 |
| car   | 5.1 | **4.9** | 2.5 |
| frog  | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**SVM Loss Example**

Georgia Tech

- If we use the softmax function to convert scores to probabilities, the right loss function to use is **cross-entropy**

- Can be derived by looking at the distance between two probability distributions (output of model and ground truth)

- Can also be derived from a maximum likelihood estimation perspective

$$s = f(x, W) \quad \textbf{Scores}$$

$$P(Y = k | X = x) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \textbf{Softmax Function}$$

$$L_i = -\log P(Y = y_i | X = x_i)$$

**Maximize log-prob of correct class = Maximize the log likelihood = Minimize the negative log likelihood**

**Performance Measure for Probabilities**

Georgia Tech

- If we use the softmax function to convert scores to probabilities, the right loss function to use is **cross-entropy**

- Goal: Minimize KL-divergence (distance measure b/w probability distributions)

$$\min_{w} KL(p^*||\hat{p}) = \sum_{y} p^*(y) \log \frac{p^*(y)}{\hat{p}(y)}$$

$$= \sum_{y} p^*(y) \log(p^*(y)) - \sum_{y} p^*(y) \log(\hat{p}(y))$$

$$p^* = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \qquad \hat{p} = \begin{bmatrix} P(Y=1|x,w) \\ P(Y=2|x,w) \\ P(Y=3|x,w) \\ P(Y=4|x,w) \\ P(Y=5|x,w) \\ P(Y=6|x,w) \\ P(Y=7|x,w) \\ P(Y=8|x,w) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.01 \\ 0.01 \\ 0.01 \\ 0.01 \\ 0.01 \\ 0.15 \\ 0.3 \end{bmatrix}$$

$-H(p^*)$
(negative entropy, term goes away because not a function of model, $W$, parameters we are minimizing over)

$H(p^*, \hat{p})$
(Cross-Entropy)

**Ground Truth**          **Prediction**

Since $p^*$ is one-hot (0 for non-ground truth classes), all we need to minimize is (where $i$ is ground truth class): $\min_{w} (-\log \hat{p}(y_i))$

**Performance Measure for Probabilities**

Georgia Tech

# **Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

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

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

**Softmax Function**

$$L_i = -\log P(Y = y_i | X = x_i)$$

Probabilities must be >= 0

Probabilities must sum to 1

| | | |
|---|---|---|
| cat | **3.2** | **24.5** |
| car | 5.1 | 164.0 |
| frog | -1.7 | 0.18 |

exp → normalize →

| **0.13** |
| 0.87 |
| 0.00 |

$$L_i = -\log(0.13)$$

**Unnormalized log-probabilities / logits**

**Unnormalized probabilities**

**Probabilities**

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**Cross-Entropy Loss Example**

Georgia Tech

# **Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

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

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

**Softmax Function**

Probabilities must be >= 0

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i | X = x_i)$$

$$L_i = -\log(0.13)$$

Q: What is the min/max of possible loss L_i?

Infimum is 0, max is unbounded (inf)

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**Cross-Entropy Loss Example**

# **Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

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

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

**Softmax Function**

Probabilities must be >= 0

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i|X = x_i)$$

$$L_i = -\log(0.13)$$

Q: At initialization all s will be approximately equal; what is the loss?

Log(C), e.g. log(10) ≈ 2

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**Cross-Entropy Loss Example**

Georgia Tech

Often, we add a **regularization term** to the loss function

**L1 Regularization**

$$L_i = |y - Wx_i|^2 + |W|$$

**Example regularizations:**

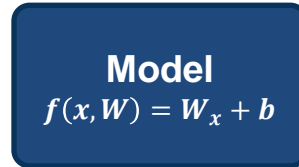- L1/L2 on weights (encourage small values)

Georgia
Tech

- Input (and representation)
- Functional form of the model
  - Including parameters
- Performance measure to improve
  - Loss or objective function
- **Algorithm for finding best parameters**
  - **Optimization algorithm**

Class Scores

Car  Coffee  Bird
     Cup

**Loss Function**

**Optimizer**

Class Scores

Car  Coffee  Bird
     Cup

**Data: Image**

**Features: Histogram**

**Model**
$$f(x, W) = W_x + b$$

Given a model and loss function, finding the best set of weights is a **search problem**

⬡ Find the best combination of weights that minimizes our loss function

**Several classes of methods:**

⬡ Random search

⬡ Genetic algorithms (population-based search)

⬡ Gradient-based optimization

In deep learning, **gradient-based methods are dominant** although not the only approach possible

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b2 \\ w_{21} & w_{22} & \cdots & w_{3m} & b3 \end{bmatrix}$$

**Loss**

# As weights change, the loss changes as well

- This is often somewhat-smooth locally, so small changes in weights produce small changes in the loss

We can therefore think about **iterative algorithms** that take **current values of weights and modify them** a bit

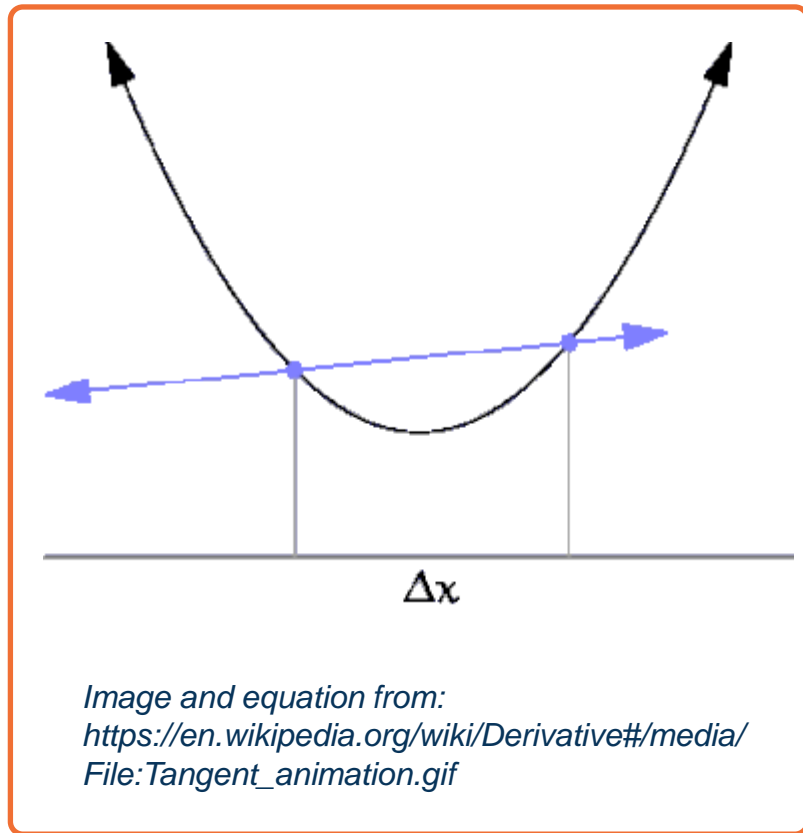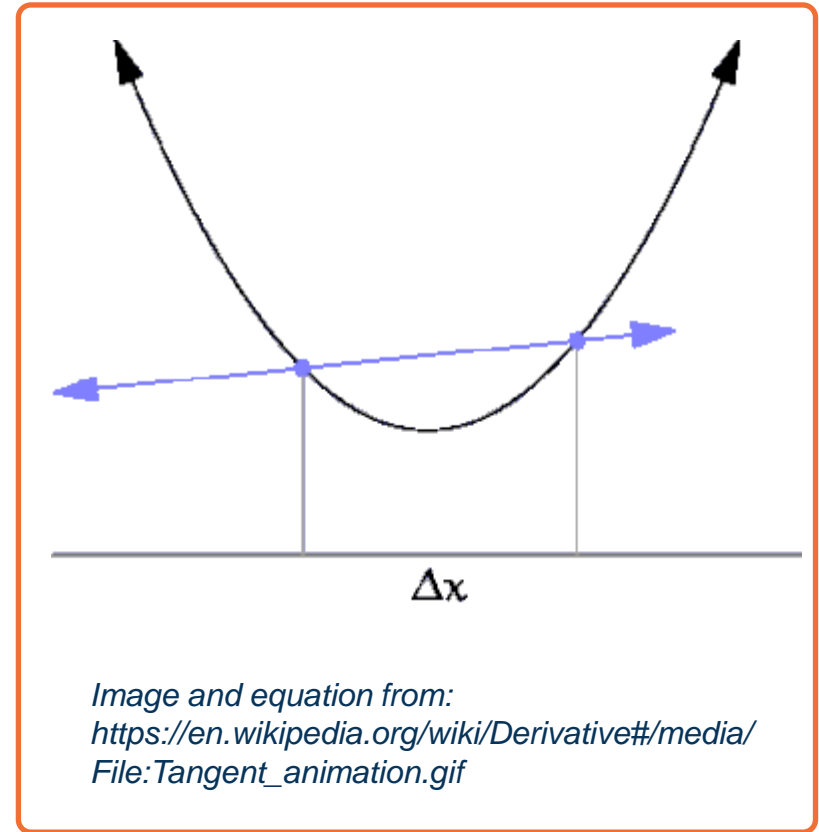Georgia Tech

Strategy: Follow the Slope!

- We can find the steepest descent direction by computing the **derivative (gradient):**

$$f'(a) = \lim_{h \to 0} \frac{f(a+h) - f(a)}{h}$$

- Steepest descent direction is the **negative gradient**

- **Intuitively:** Measures how the function changes as the argument a changes by a small step size

  - As step size goes to zero

- **In Machine Learning:** Want to know how the **loss function** changes **as weights** are varied

  - Can consider each parameter separately by taking **partial derivative** of loss function with respect to that parameter



*Image and equation from: https://en.wikipedia.org/wiki/Derivative#/media/File:Tangent_animation.gif*

**Derivatives**

Georgia Tech

$$f'(a) = \lim_{h \to 0} \frac{f(a+h) - f(a)}{h}$$



*Image and equation from:*
*https://en.wikipedia.org/wiki/Derivative#/media/*
*File:Tangent_animation.gif*

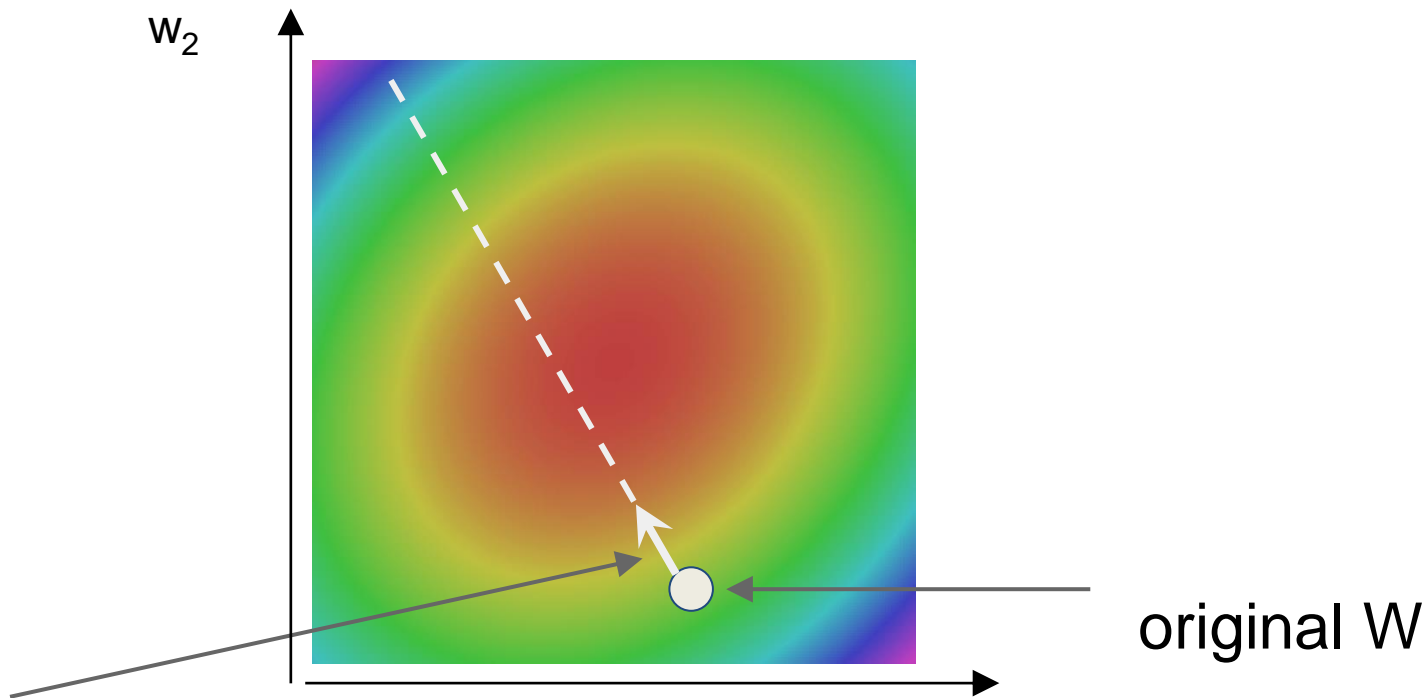**Derivatives in d-dimensions**

Georgia Tech

This idea can be turned into an **algorithm (gradient descent)**

1. Choose a model: $f(x, W) = \mathrm{W}x$

2. Choose loss function: $L_i = (y - Wx_i)^2$

3. Calculate partial derivative for each parameter: $\dfrac{\partial L}{\partial w_i}$

4. Update the parameters: $w_i = w_i - \dfrac{\partial L}{\partial w_i}$

   Add learning rate to prevent too big of a step: $w_i = w_i - \alpha \dfrac{\partial L}{\partial w_i}$
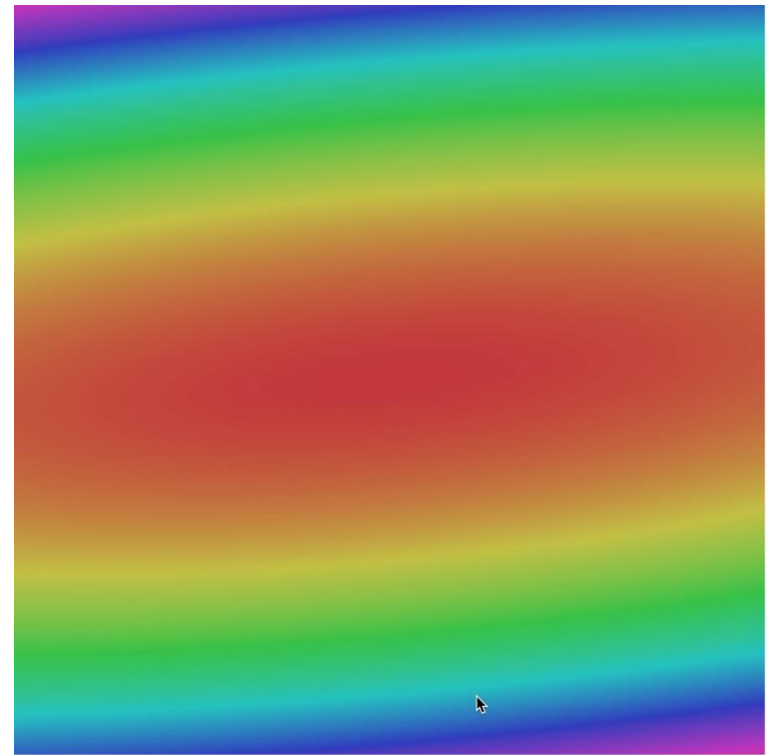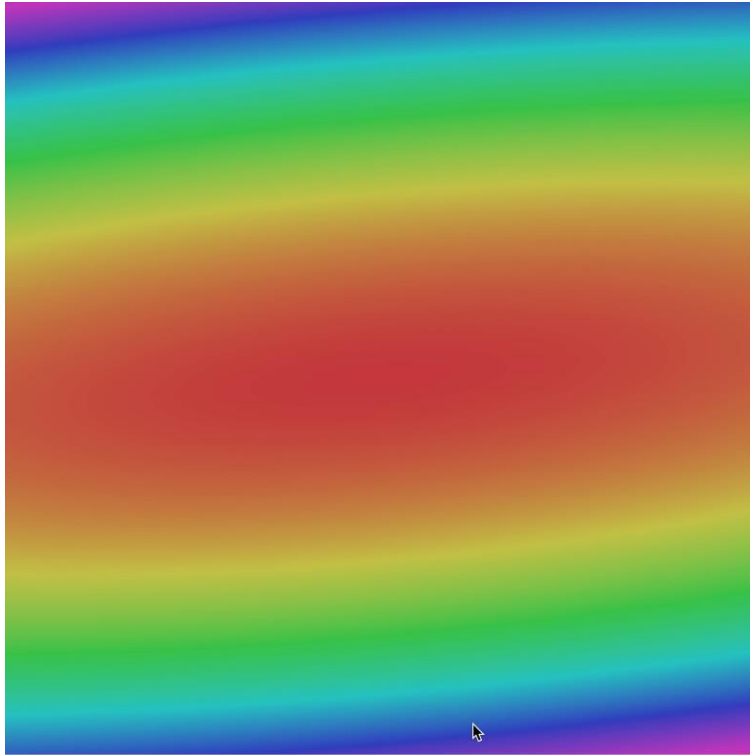
5. Repeat (from Step 3)

**Gradient Descent**

Georgia Tech

$w_2$

original W

negative gradient direction

$w_1$

**Gradient Descent**

Georgia Tech

$$w_1$$

**Gradient Descent**

Often, we only compute the gradients across a small subset of data

- Full Batch Gradient Descent

$$L = \frac{1}{N} \sum L\left(f(x_i, W), y_i\right)$$

- Mini-Batch Gradient Descent

$$L = \frac{1}{M} \sum L\left(f(x_i, W), y_i\right)$$

   - Where M is a *subset* of data

- We iterate over mini-batches:

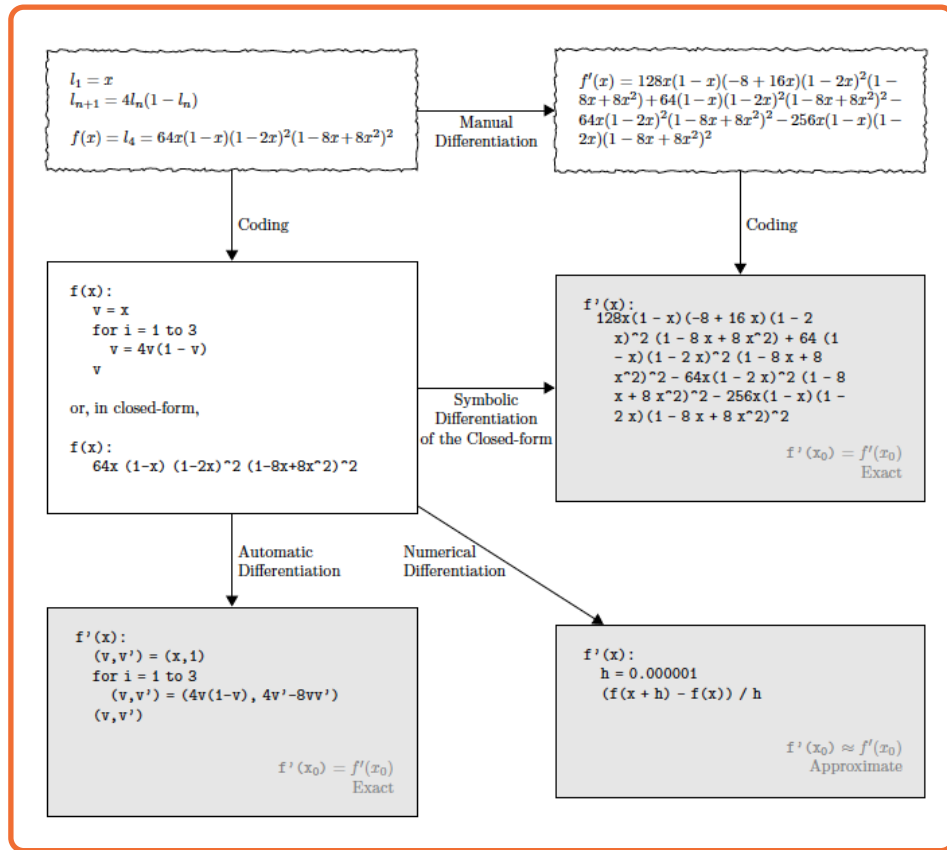   - Get mini-batch, compute loss, compute derivatives, and take a set

Gradient descent is guaranteed to converge under some conditions

- For example, learning rate has to be appropriately reduced throughout training

- It will converge to a *local* minima

    - Small changes in weights would not decrease the loss

- It turns out that some of the local minima that it finds in practice (if trained well) are still pretty good!

Georgia Tech

We know how to compute the **model output and loss function**

**Several ways to compute** $\frac{\partial L}{\partial w_i}$

- ⬡ Manual differentiation

- ⬡ Symbolic differentiation

- ⬡ Numerical differentiation

- ⬡ Automatic differentiation

Georgia Tech

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**gradient dW:**

[?,
?,
?,
?,
?,
?,
?,
?,
?,…]

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (first dim)**:**

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25322**

**gradient dW:**

[?,
?,
?,
?,
?,
?,
?,
?,
?,…]

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (first dim)**:**

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25322**

**gradient dW:**

[**-2.5**,
?,
?,

(1.25322 - 1.25347)/0.0001
= -2.5

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

?,
?,…]

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (second dim)**:**

[0.34,
-1.11 **+ 0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25353**

**gradient dW:**

[-2.5,
?,
?,
?,
?,
?,
?,
?,
?,…]

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (second dim)**:**

[0.34,
-1.11 **+ 0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25353**

**gradient dW:**

[-2.5,
**0.6**,
?,
?,

(1.25353 - 1.25347)/0.0001
= 0.6

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

?,…]

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**W + h** (third dim)**:**

[0.34,
-1.11,
0.78 + **0.0001**,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**gradient dW:**

[-2.5,
0.6,
?,
?,
?,
?,
?,
?,…]

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**
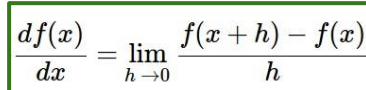
**W + h** (third dim)**:**

[0.34,
-1.11,
0.78 + **0.0001**,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
**loss 1.25347**

**gradient dW:**

[-2.5,
0.6,
**0**,
?,

(1.25347 - 1.25347)/0.0001
= 0

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

?,…]

# Numerical vs Analytic Gradients

$$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

**Numerical gradient**: slow :(, approximate :(, easy to write :)
**Analytic gradient**: fast :), exact :), error-prone :(

In practice: Derive analytic gradient, check your implementation with numerical gradient.
This is called a **gradient check.**

- Components of parametric classifiers:
  - Input/Output: Image/Label
  - Model (function): Linear Classifier + Softmax
  - Loss function: Cross-Entropy
  - Optimizer: Gradient Descent

- Ways to compute gradients
  - Numerical
  - Next: Analytical, automatic differentiation

Georgia Tech

# For some functions, we can analytically derive the partial derivative

## Example:

### Function

$$f(w, x_i) = w^T x_i$$

(Assume $w$ and $x_i$ are column vectors, so same as $w \cdot x_i$)

### Loss

$$(y_i - w^T x_i)^2$$

### Update Rule

$$w_j \leftarrow w_j + 2\alpha \sum_{k=1}^{N} \delta_k x_{kj}$$

## Derivation of Update Rule

Georgia Tech

# For some functions, we can analytically derive the partial derivative

## Example:

### Function

$$f(w, x_i) = w^T x_i$$

(Assume $w$ and $x_i$ are column vectors, so same as $w \cdot x_i$)

### Loss

$$(y_i - w^T x_i)^2$$

### Update Rule

$$w_j \leftarrow w_j + 2\alpha \sum_{k=1}^{N} \delta_k x_{kj}$$

## Derivation of Update Rule

$$L = \sum_{k=1}^{N} (y_k - w^T x_k)^2$$

$$\frac{\partial L}{\partial w_j} = \sum_{k=1}^{N} \frac{\partial}{\partial w_j} (y_k - w^T x_k)^2$$

Gradient descent tells us we should update $w$ as follows to minimize $L$:

$$= \sum_{k=1}^{N} 2(y_k - w^T x_k) \frac{\partial}{\partial w_j} (y_k - w^T x_k)$$

$$= -2 \sum_{k=1}^{N} \delta_k \frac{\partial}{\partial w_j} w^T x_k$$

$$w_j \leftarrow w_j - \eta \frac{\partial L}{\partial w_j}$$

…where…
$$\delta_k = y_k - w^T x_k$$

So what's $\frac{\partial L}{\partial w_j}$?

$$= -2 \sum_{k=1}^{N} \delta_k \frac{\partial}{\partial w_j} \sum_{i=1}^{m} w_i x_{ki}$$

$$= -2 \sum_{k=1}^{N} \delta_k x_{kj}$$

## Manual Differentiation

Georgia Tech

# If we add a **non-linearity (sigmoid),** derivation is more complex

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$
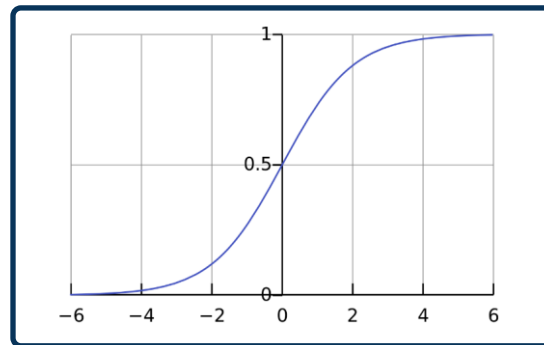
First, one can derive that: $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

$$\mathbf{f(x)} = \sigma\left(\sum_k w_k x_k\right)$$

$$L = \sum_i \left(y_i - \sigma\left(\sum_k w_k x_{ik}\right)\right)^2$$

$$\frac{\partial L}{\partial w_j} = \sum_i 2\left(y_i - \sigma\left(\sum_k w_k x_{ik}\right)\right)\left(-\frac{\partial}{\partial w_j}\sigma\left(\sum_k w_k x_{ik}\right)\right)$$

$$= \sum_i -2\left(y_i - \sigma\left(\sum_k w_k x_{ik}\right)\right)\sigma'\left(\sum_k w_k x_{ik}\right)\frac{\partial}{\partial w_j}\sum_k w_k x_{ik}$$

$$= \sum_i -2\delta_i\sigma(\mathbf{d_i})(1 - \sigma(\mathbf{d_i}))x_{ij}$$

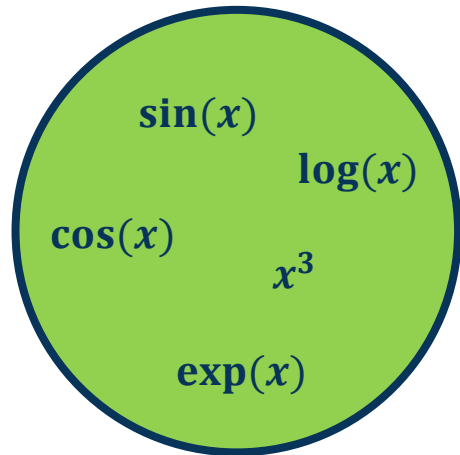where $\delta_i = y_i - \mathbf{f(x_i)}$ $\qquad d_i = \sum w_k x_{ik}$

## The sigmoid perception update rule:

$$w_j \leftarrow w_j + 2\alpha\sum_{k=1}^N \delta_i\sigma_i(1 - \sigma_i)x_{ij}$$

where $\sigma_i = \sigma\left(\sum_{j=1}^m w_j x_{ij}\right)$

$$\delta_i = y_i - \sigma_i$$

**Adding a Non-Linear Function**

Georgia Tech

Given a library of simple functions



$\sin(x)$

$\log(x)$

$\cos(x)$

$x^3$

$\exp(x)$

Compose into a complicate function

$$-\log\left(\frac{1}{1 + e^{-w \cdot x}}\right)$$

$$w \cdot x \xrightarrow{u} \frac{1}{1 + e^{-u}} \xrightarrow{p} -\log(p) \xrightarrow{L}$$

*Adapted from slides by: Marc'Aurelio Ranzato, Yann LeCun*

**Decomposing a Function**

Georgia Tech