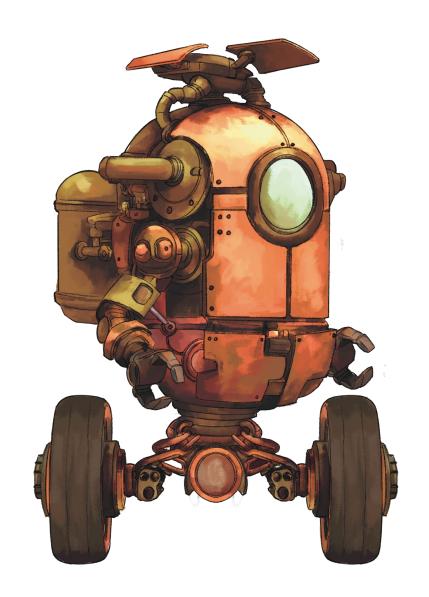
CS 3630, Fall 2025

Lecture 19: Learning CNNs



Many slides adapted from Stanford's CS231N by Fei-Fei Li, Justin Johnson, Serena Yeung, as well as Slides by Marc'Aurelio Ranzato (NYU), Dhruv Batra & Devi Parikh (Georgia Tech)

Outline

- 1. Intra-class variability
- 2. Supervised Learning
- 3. Regression and Classification Losses
- 4. Stochastic Gradient Descent
- 5. Calculating Gradients

Image Classification: A core task in Computer Vision

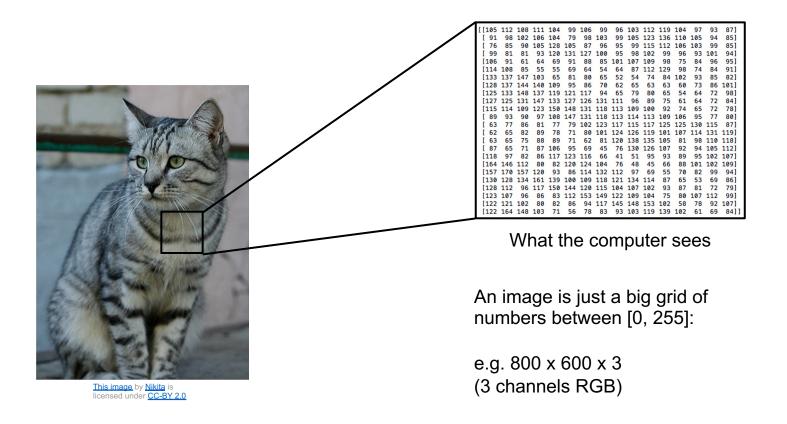


This image by Nikita is licensed under CC-BY 2.0

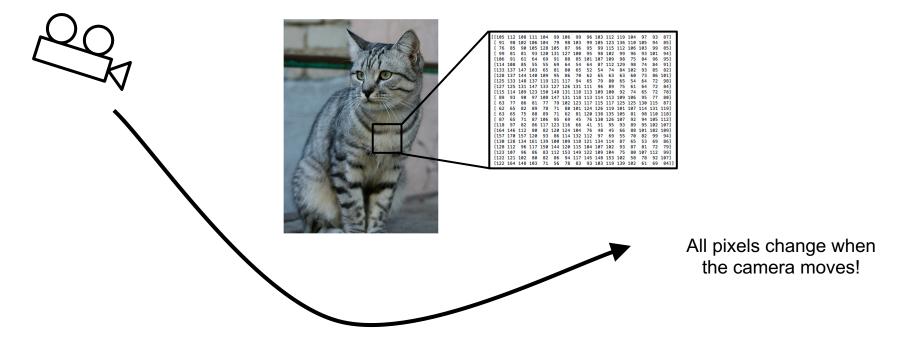
(assume given set of discrete labels) {dog, cat, truck, plane, ...}

----- cat

The Problem: Semantic Gap



Challenges: Viewpoint variation



This image by Nikita is licensed under CC-BY 2.0

Challenges: Illumination



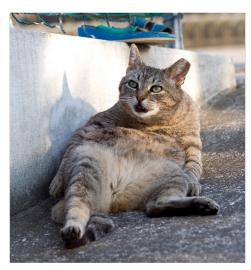






This image is CC0 1.0 public domain

Challenges: Deformation



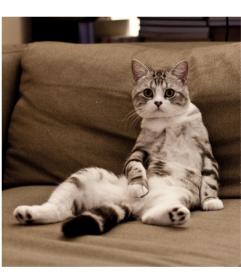
This image by Umberto Salvagnin is licensed under CC-BY 2.0



This image by Umberto Salvagnin is licensed under CC-BY 2.0



<u>This image</u> by <u>sare bear</u> is licensed under <u>CC-BY 2.0</u>



This image by Tom Thai is licensed under CC-BY 2.0

Challenges: Occlusion







This image is CC0 1.0 public domain

This image is CC0 1.0 public domain

This image by ionsson is licensed under CC-BY 2.0

Challenges: Background Clutter





This image is CC0 1.0 public domain

This image is CC0 1.0 public domain

Challenges: Intraclass variation



This image is CC0 1.0 public domain

Outline

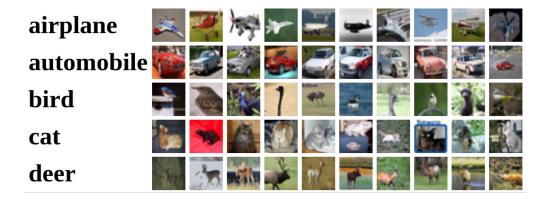
- 1. Intra-class variability
- 2. Supervised Learning
- 3. Regression and Classification Losses
- 4. Stochastic Gradient Descent
- 5. Calculating Gradients

ML: A Data-Driven Approach

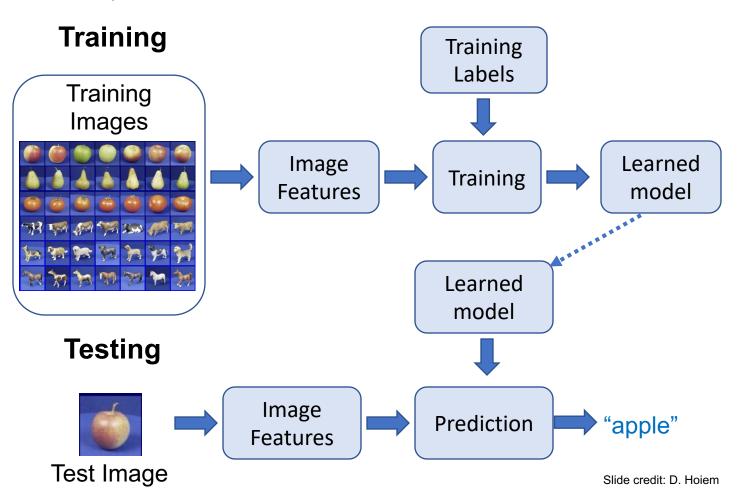
- 1. Collect a dataset of images x and labels y
- 2. Use Machine Learning to train a classifier
- 3. Evaluate the classifier on new images

def train(images, labels): # Machine learning! return model def predict(model, test_images): # Use model to predict labels return test_labels

Example training set



Steps



Outline

- 1. Intra-class variability
- 2. Supervised Learning
- 3. Regression and Classification Losses
- 4. Stochastic Gradient Descent
- 5. Calculating Gradients

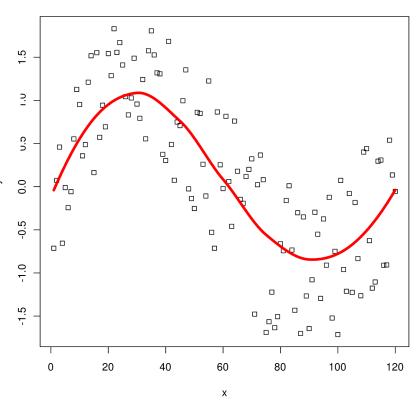
Two different learning problems are classification and regression

- Regression: continuous labels y
- Sum-square-difference loss:

$$L_{ ext{SSD}}(W;D) \doteq \sum_{(x,y) \in D} |f(x;W) - y|^2$$

• Example: direct 3D pose regression:





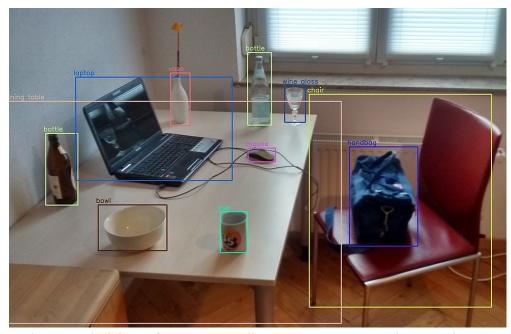
By Kierano - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=7034448

Two different learning problems are classification and regression

- Classification: discrete labels y
- Cross-entropy loss:

$$\mathcal{L}_{ ext{CE}}(heta;D) \doteq \sum_{c} \sum_{(x,y=c) \in D} \log rac{1}{p_c(x; heta)}$$

- Average surprise!
- Example: object detection:



(MTheiler), CC BY-SA 4.0 https://creativecommons.org/licenses/by-sa/4.0, via Wikimedia Commons

How to get probabilities? **Softmax** converts a set of C class "scores" into "probabilities"

- If we need to classify inputs into C different classes, we put C units in the last layer to produce C *one-vs.-others* scores $f_1, f_2, ..., f_C$
- Apply *softmax* function to convert these scores to probabilities:

$$\operatorname{softmax}(f_1, \dots, f_c) = \left(\frac{\exp(f_1)}{\sum_j \exp(f_j)}, \dots, \frac{\exp(f_C)}{\sum_j \exp(f_j)}\right)$$

If one of the inputs is much larger than the others, then the corresponding softmax value will be close to 1 and others will be close to 0

Outline

- 1. Intra-class variability
- 2. Supervised Learning
- 3. Regression and Classification Losses
- 4. Stochastic Gradient Descent
- 5. Calculating Gradients

How to minimize the loss by changing the weights? Strategy: Follow the slope of the loss function



Strategy: Follow the slope

In 1-dimension, the derivative of a function:

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

In multiple dimensions, the gradient is the vector of (partial derivatives) along each dimension

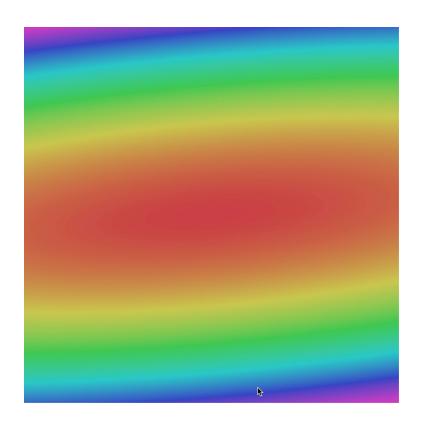
The slope in any direction is the **dot product** of the direction with the gradient

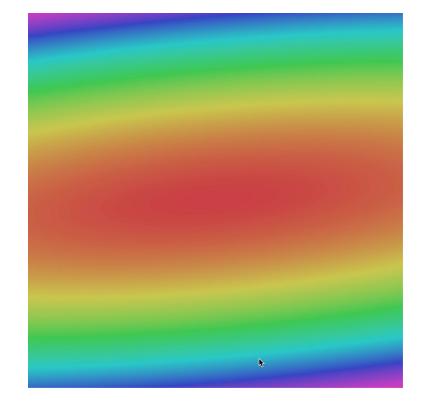
The direction of steepest descent is the negative gradient

Gradient Descent

```
# Vanilla Gradient Descent

while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```





Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(x_i, y_i, W) + \lambda R(W)$$

$$1 \sum_{i=1}^{N} L_i(x_i, y_i, W) + \lambda R(W)$$

 $\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$

Full sum expensive when N is large!

Approximate sum using a **minibatch** of examples 32 / 64 / 128 common

Vanilla Minibatch Gradient Descent

while True:

```
data_batch = sample_training_data(data, 256) # sample 256 examples
weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
weights += - step_size * weights_grad # perform parameter update
```

Outline

- 1. Intra-class variability
- 2. Supervised Learning
- 3. Regression and Classification Losses
- 4. Stochastic Gradient Descent
- 5. Calculating Gradients

How do we *really* compute gradients?

- Analytic or "Manual" Differentiation
- Symbolic Differentiation
- Numerical Differentiation
- Automatic Differentiation!
 - Forward mode AD
 - Reverse mode AD
 - aka "backpropagation"
 - Implemented in specialized frameworks:
 - pytorch (Facebook)
 - TensorFlow (Google) frameworks
 - Main computation, mainly done on GPU (or TPU)



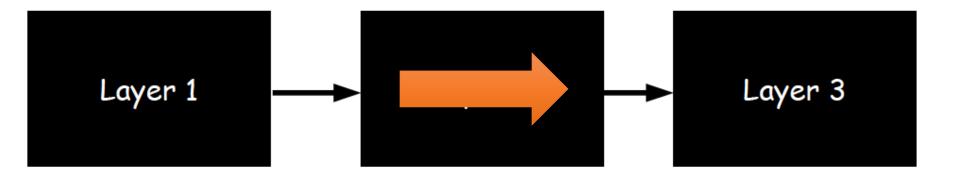
• Step 1: Compute Loss on mini-batch

[F-Pass]



• Step 1: Compute Loss on mini-batch

[F-Pass]



• Step 1: Compute Loss on mini-batch

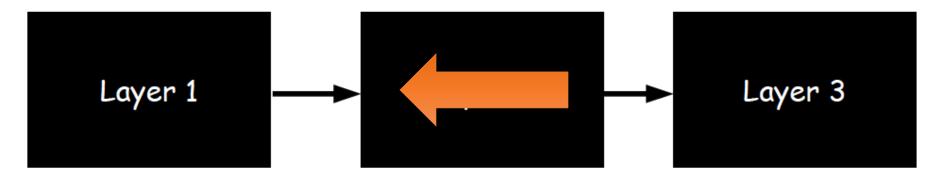
[F-Pass]



- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients wrt parameters [B-Pass]



- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients wrt parameters [B-Pass]



- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients wrt parameters [B-Pass]



Step 1: Compute Loss on mini-batch

- [F-Pass]
- Step 2: Compute gradients wrt parameters
- [B-Pass]

Step 3: Use gradient to update parameters



$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

Outline

- 1. Intra-class variability: viewpoint, lighting, instance
- 2. Supervised Learning: label + optimization
- 3. Regression and Classification: SSD + CE
- 4. Stochastic Gradient Descent: mini-batches
- 5. Calculating Gradients: back-propagation