

Topics:

- Machine learning intro, applications (CV, NLP, etc.)
- Parametric models and their components

CS 4644 / 7643-A
ZSOLT KIRA

- **PS0 due May 18th Sunday night, but do it TODAY!**
 - Please do it, and give others a chance at waitlist if your background is not sufficient (beef it up and take it next time)
 - Do it even if you're on the waitlist!
- **Piazza:**
 - <https://piazza.com/gatech/summer2025/cs46447643a/>
 - Search for teammates: @5 (<https://piazza.com/class/mafsg3dobtu42c/post/5>)
 - Note: Do NOT post anything containing solutions publicly!
 - Make it active!
- **Office hours** start next week

- **Collaboration**
 - Only on HWs and project (not allowed in HW0/PS0).
 - You may discuss the questions
 - Each student writes their own answers
 - Write on your homework anyone with whom you collaborate
 - Each student must write their own code for the programming part
 - Do NOT search for code implementing what we ask; search for concepts
- **Zero tolerance on plagiarism**
 - Neither ethical nor in your best interest
 - Always credit your sources
 - Don't cheat. We will find out.

- **Grace period**
 - 2 days grace period for each assignment (**EXCEPT PS0**)
 - Intended for checking submission NOT to replace due date
 - No need to ask for grace, no penalty for turning it in within grace period
 - Can NOT use for PS0
- **After grace period, you get a 0 (no excuses except medical)**
 - Send all medical requests to dean of students (<https://studentlife.gatech.edu/>)
 - Form: https://gatech-advocate.symplicity.com/care_report/index.php/pid224342
- **DO NOT SEND US ANY MEDICAL INFORMATION!** We do not need any details, just a confirmation from dean of students

Python Numpy Tutorial

This tutorial was contributed by [Justin Johnson](#).

We will use the Python programming language for all assignments in this course. Python is a great general-purpose programming language on its own, but with the help of a few popular libraries (numpy, scipy, matplotlib) it becomes a powerful environment for scientific computing.

We expect that many of you will have some experience with Python and numpy; for the rest of you, this section will serve as a quick crash course both on the Python programming language and on the use of Python for scientific computing.

<http://cs231n.github.io/python-numpy-tutorial/>

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

Python Numpy Tutorial

This tutorial was contributed by [Justin Johnson](#).

We will use the Python programming language for all assignments in this course. Python is a great general-purpose programming language on its own, but with the help of a few popular libraries (numpy, scipy, matplotlib) it becomes a powerful environment for scientific computing.

We expect that many of you will have some experience with Python and numpy; for the rest of you, this section will serve as a quick crash course both on the Python programming language and on the use of Python for scientific computing.

<http://cs231n.github.io/python-numpy-tutorial/>

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

What is Machine Learning (ML)?

“A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .”

Tom Mitchell (Machine Learning, 1997)

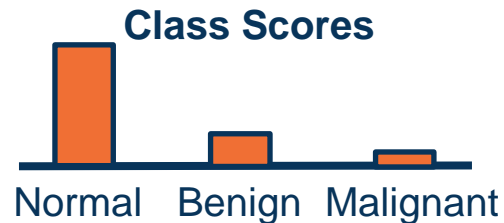
Application: Computer Vision



Model



Model



Model



3D Reconstructions



Kanazawa et al., 2018

Example: Image Classification

Supervised Learning

- Train Input: $\{X, Y\}$
- Learning output: $f : X \rightarrow Y$,
e.g. $P(y|x)$

Unsupervised Learning

- Input: $\{X\}$
- Learning output: $P(x)$
- Example: Clustering, density estimation, etc.

Reinforcement Learning

- Supervision in form of **reward**
- No supervision on what action to take

Very often combined, sometimes within the same model!

Parametric Model

Explicitly model the function $f : X \rightarrow Y$ in the form of a parametrized function $f(x, W) = y$, **examples**:

- ⬡ Logistic regression/classification
- ⬡ Neural networks

Capacity (size of hypothesis class) **does not** grow with size of training data!

Learning is **search**

Parametric – Linear Classifier

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

Training Stage:

Training Data $\{ (x_i, y_i) \} \rightarrow h$ (Learning)

Testing Stage

Test Data $x \rightarrow h(x)$ (Apply function, Evaluate error)

Probabilities to rescue:

X and Y are *random variables*

$$D = (x_1, y_1), (x_2, y_2), \dots, (x_N, y_N) \sim P(X, Y)$$

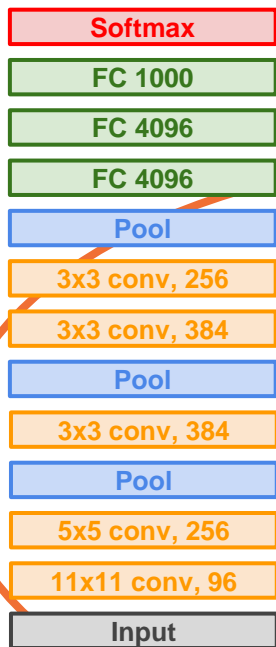
IID: Independent Identically Distributed

Both training & testing data sampled IID from $P(X, Y)$

Learn on training set

Have some hope of *generalizing* to test set

AlexNet



model class

Optimization Error (Poor/simple optimizer)

Estimation Error (Finite data)

Modeling Error (gap from reality)

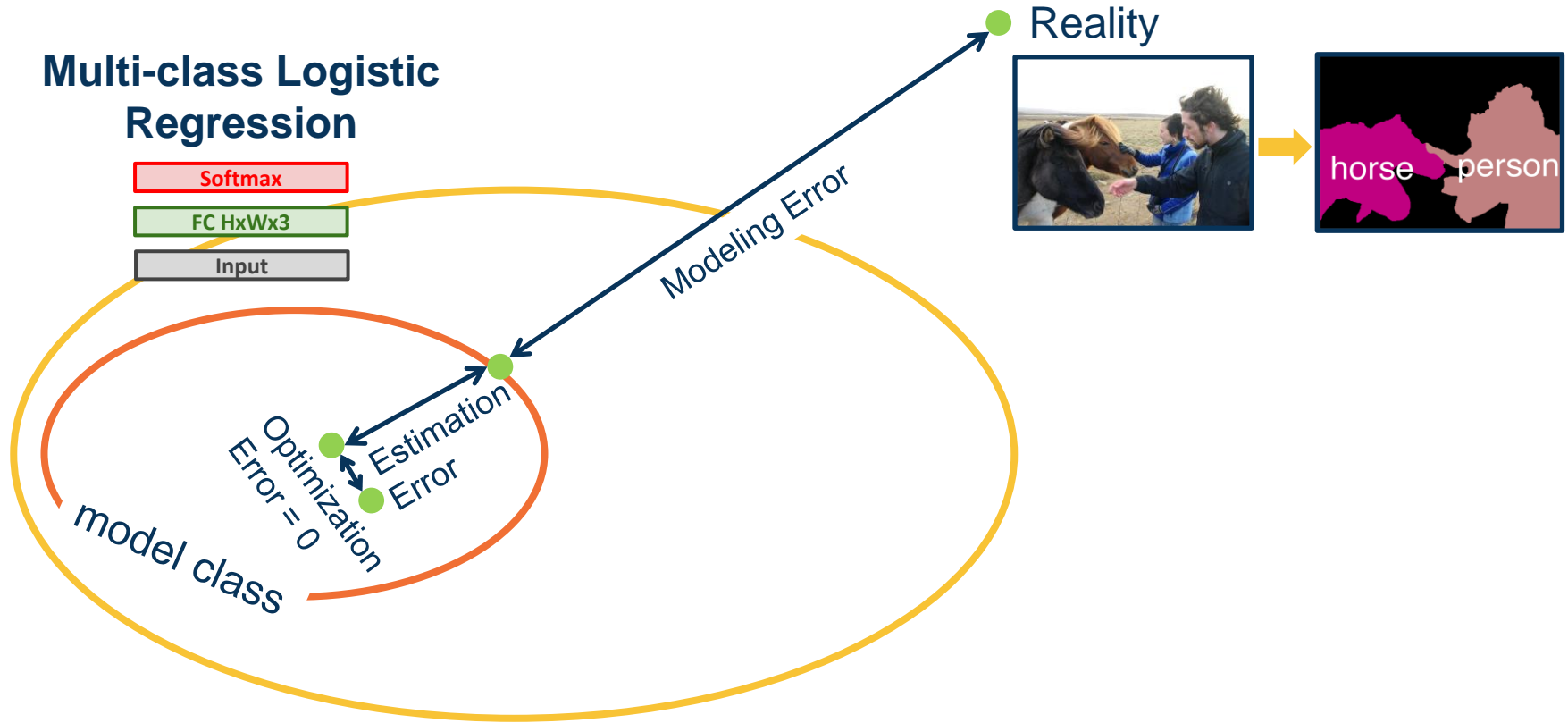
Reality



From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Generalization

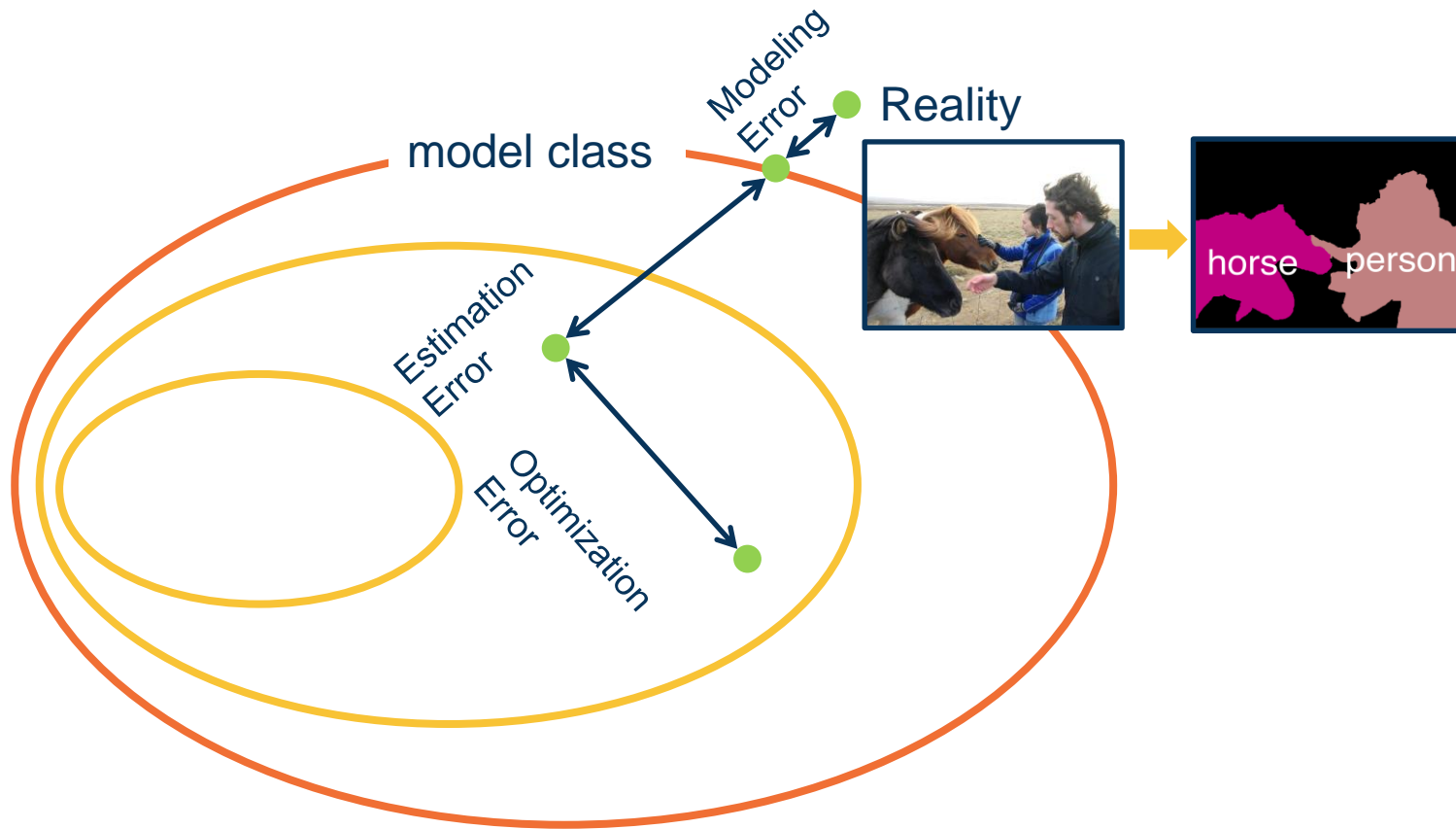
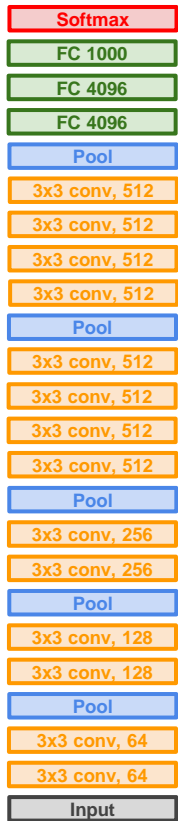
Multi-class Logistic Regression



From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Generalization

VGG19



From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Generalization

20 years of research in Learning Theory oversimplified:

If you have:

Enough training data D

and H is not too complex

then *probably* we can generalize to unseen test data

Caveats: A number of recent empirical results question our intuitions built from this clean separation.

Zhang et al., Understanding deep learning requires rethinking generalization

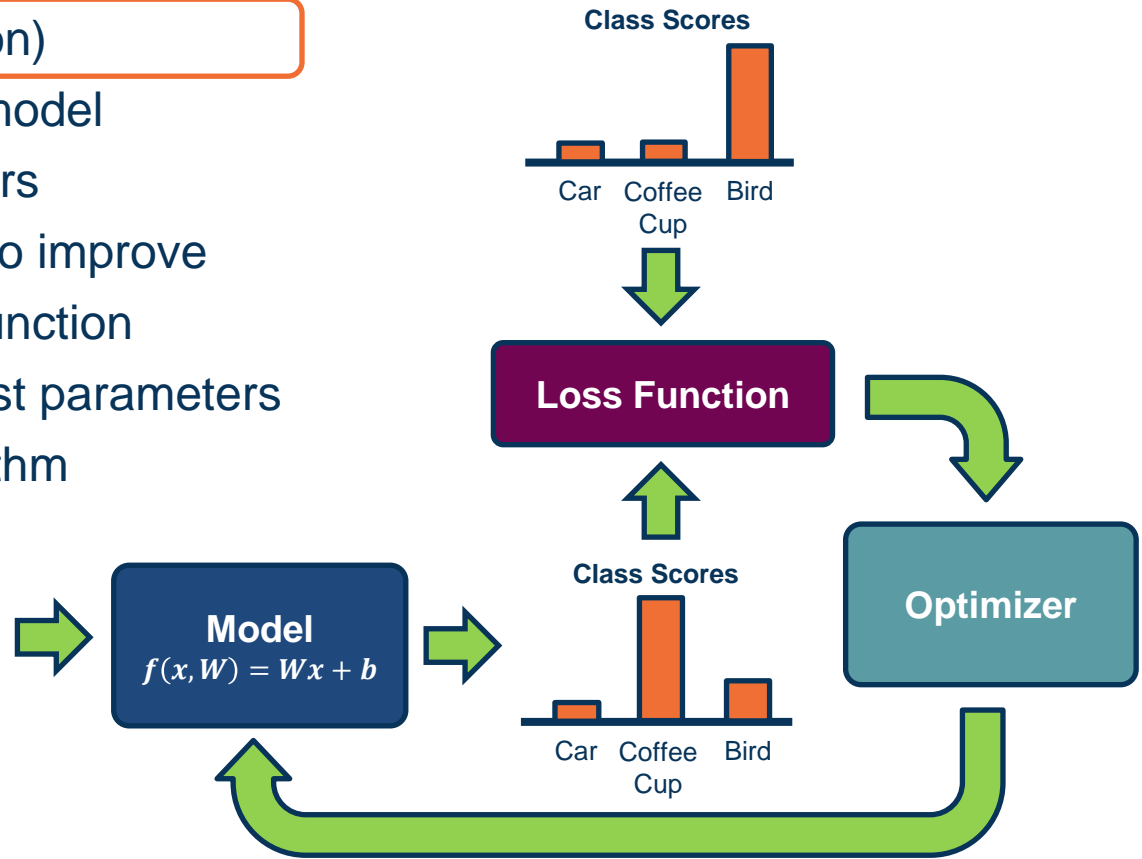
Guarantees

Components of a Parametric Learning Algorithm

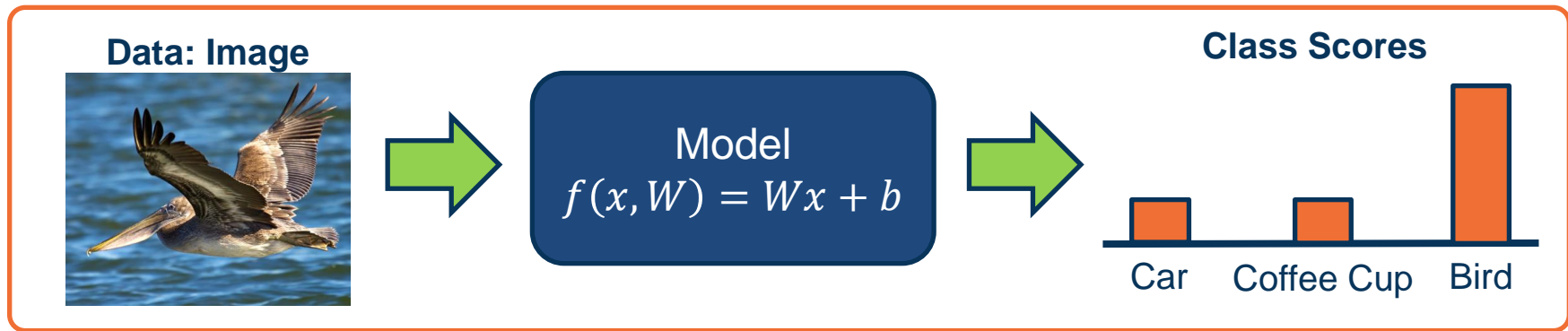
- Input (and representation)
- Functional form of the model
 - Including parameters
- Performance measure to improve
 - Loss or objective function
- Algorithm for finding best parameters
 - Optimization algorithm



Data: Image



Components of a Parametric Model



Input $\{X, Y\}$ where:

- ✦ X is an image
- ✦ Y is a **ground truth label** annotated by an expert (human)
- ✦ $f(x, W) = Wx + b$ is our model, chosen to be a linear function in this case
- ✦ W and b are the parameters (**weights**) of our model that must be learned

Example: Image Classification

Input image is **high-dimensional**

- For example $n=512$ so 512×512 image = **262,144** pixels
- Learning a classifier with high-dimensional inputs is hard

Before deep learning, it was typical to perform **feature engineering**

- Hand-design algorithms for converting raw input into a lower-dimensional set of features

Input Image



$$x = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix}$$

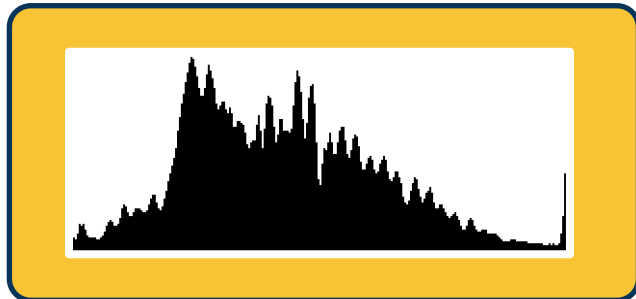
Example: Color histogram

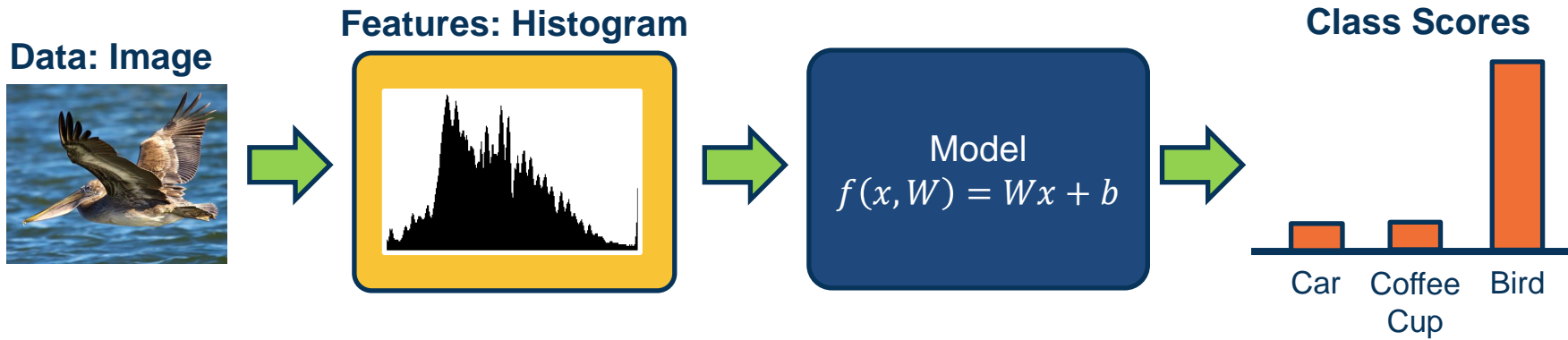
- Vector of numbers representing number of pixels fitting within each bin
- We will later see that learning the feature representation itself is much more effective

Data: Image



Features: Histogram



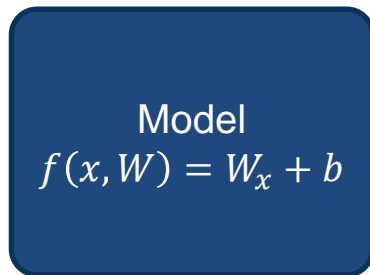


Input $\{X, Y\}$ where:

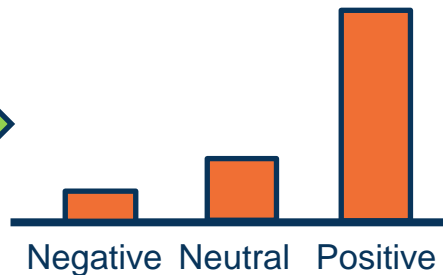
- ✦ X is an **image histogram**
- ✦ Y is a **ground truth label represented a probability distribution**
- ✦ $f(x, W) = Wx + b$ is our model, chosen to be a linear function in this case
- ✦ W and b are the **weights** of our model that must be learned

Example: Image Classification

Data: Text



Class Scores



Input $\{X, Y\}$ where:

- X is a sentence
- Y is a **ground truth label** annotated by an expert (human)
- $f(x, W) = Wx + b$ is our model, chosen to be a linear function in this case
- W and b are the **weights** of our model that must be learned

Word Histogram

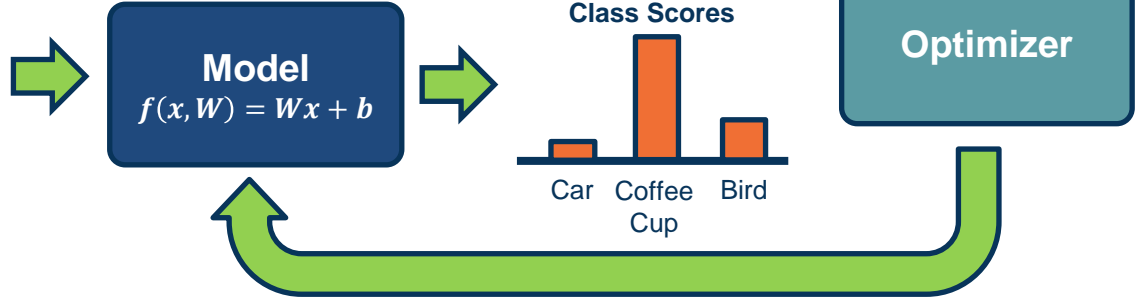
Word	Count
this	1
that	0
is	2
...	
extremely	1
hello	0
onomatopoeia	0
...	

Example: Image Classification

- 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



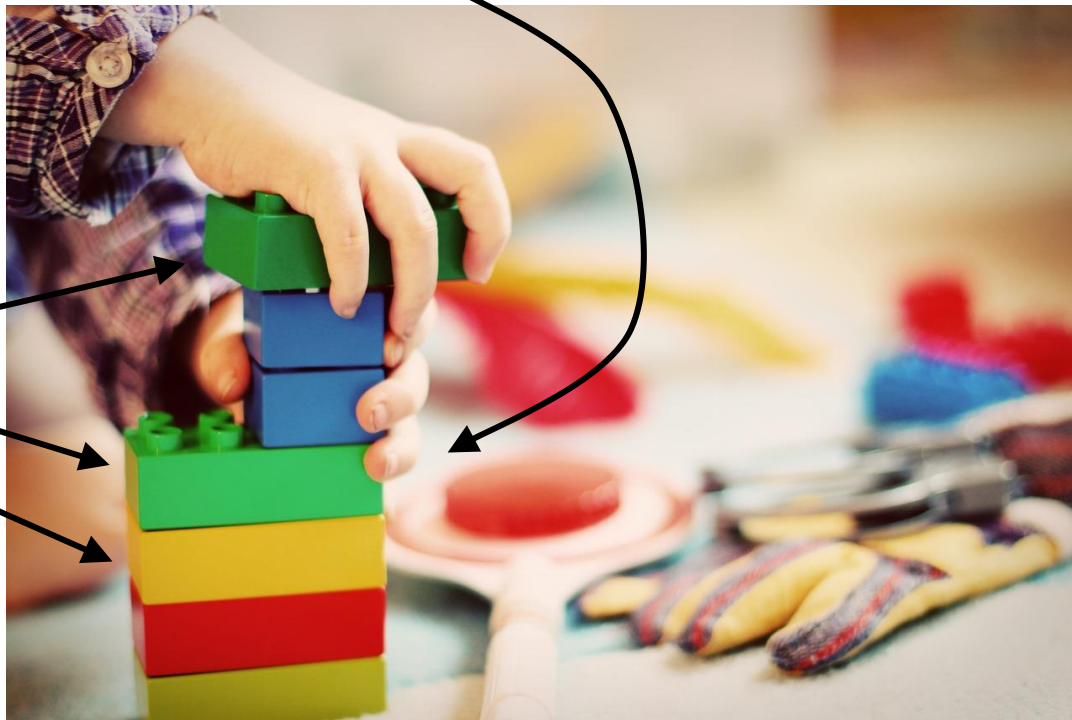
Data: Image



Components of a Parametric Model

Neural Network

Linear
classifiers

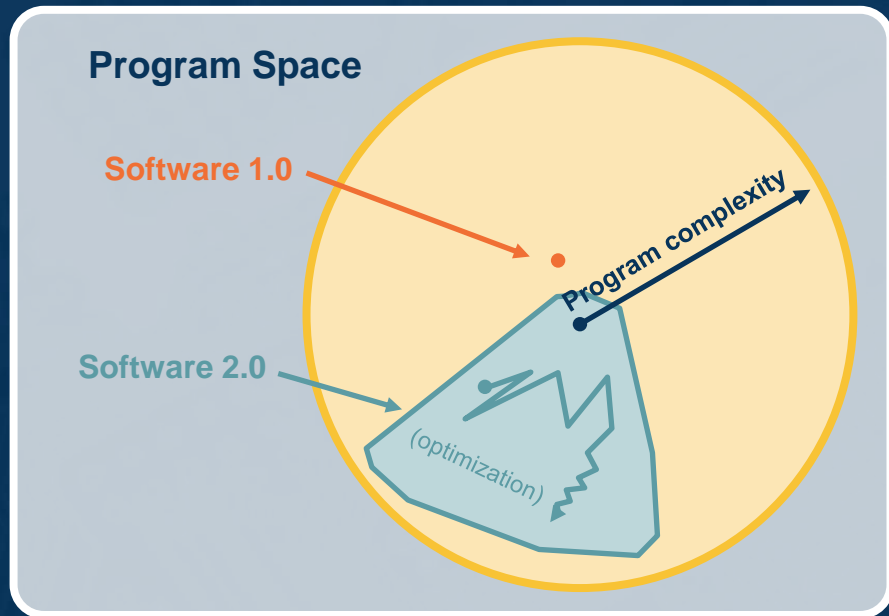
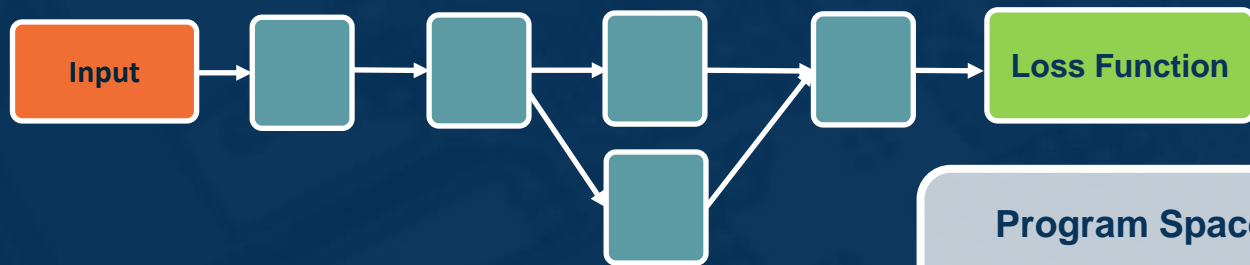


[This image](#) is [CC0 1.0](#) public domain

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

Deep Learning as Legos

The Power of Deep Learning

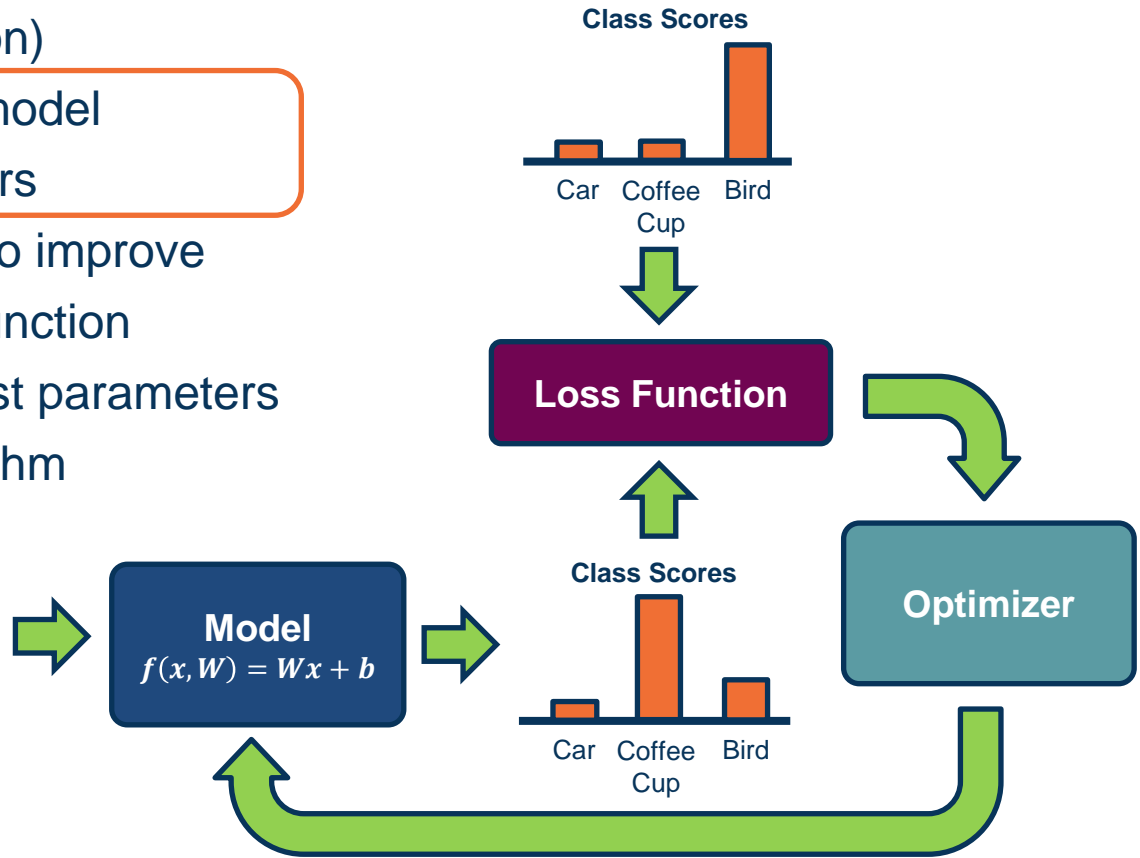


Adapted from figure by Andrej Karpathy

- 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



Components of a Parametric Model

What is the **simplest** function
you can think of?

Car



Bird



Our model is:

$$f(x, w) = w \cdot x + b$$

**Classifier
Result**

Weights

Input

**Bias
(scalar)**

(Note if w and x are column vectors we often show this as $w^T x$)

Simple Function

Linear Classification and Regression

Simple linear classifier:

- Calculate score:

$$f(x, w) = w \cdot x + b$$

- Binary classification rule (w is a vector):

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

- For multi-class classifier take class with highest (max) score

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



Data: Image



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



Class Scores



$$x = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix} \xrightarrow{\text{Flatten}} x = \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{21} \\ x_{22} \\ \vdots \\ x_{n1} \\ \vdots \\ x_{nn} \end{bmatrix}$$

To simplify notation we will refer to inputs as $x_1 \cdots x_m$ where $m = n \times n$

Input Dimensionality

$$\text{Model}$$

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

Classifier for class 1	→	W_{11}	W_{12}	\cdots	W_{1m}	$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$
Classifier for class 2	→	W_{21}	W_{22}	\cdots	W_{2m}	
Classifier for class 3	→	W_{31}	W_{32}	\cdots	W_{3m}	

W
 x
 b

(Note that in practice, implementations can use xW instead, assuming a different shape for W . That is just a different convention and is equivalent.)

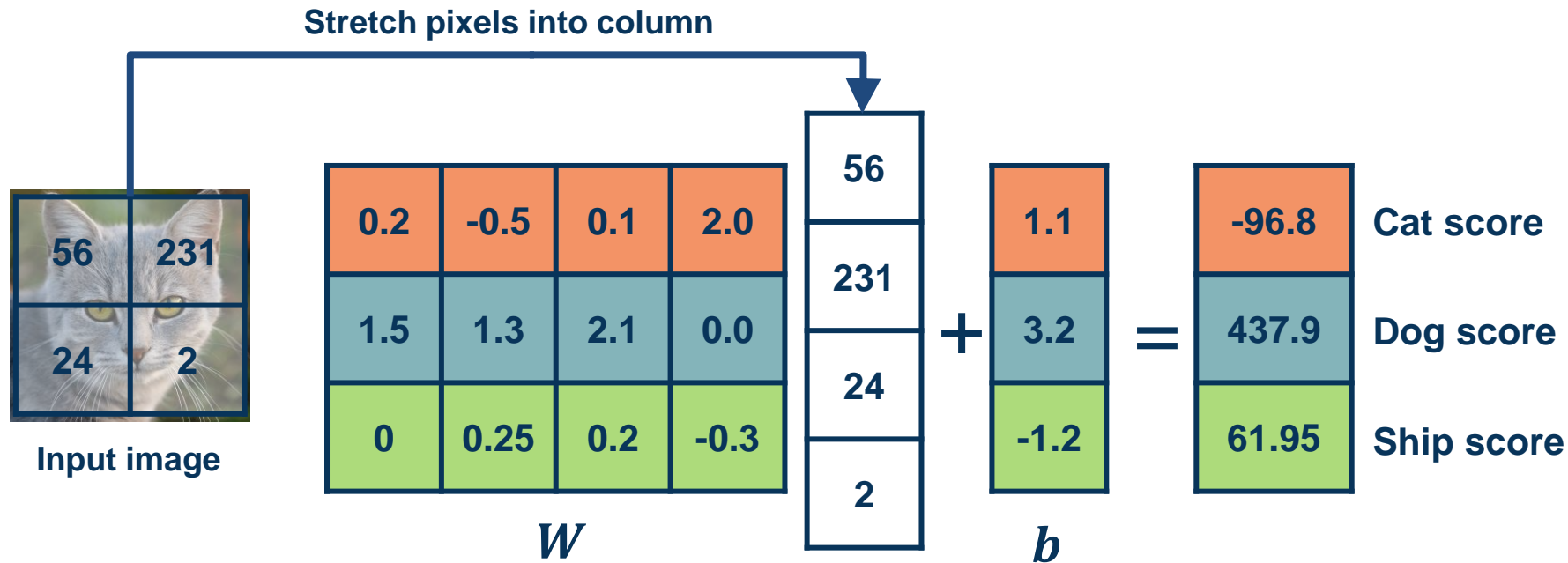
- We can move the bias term into the weight matrix, and a “1” at the end of the input
- Results in **one matrix-vector multiplication!**

$$\text{Model}$$
$$f(x, W) = Wx + b$$

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

W x

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



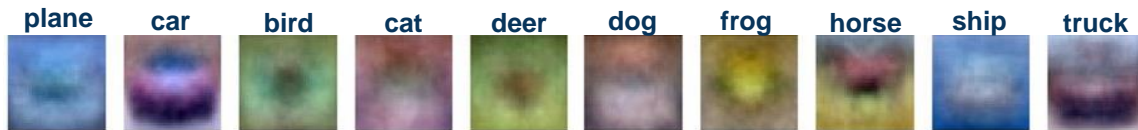
Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck



Visual Viewpoint

We can convert the weight vector back into the shape of the image and visualize



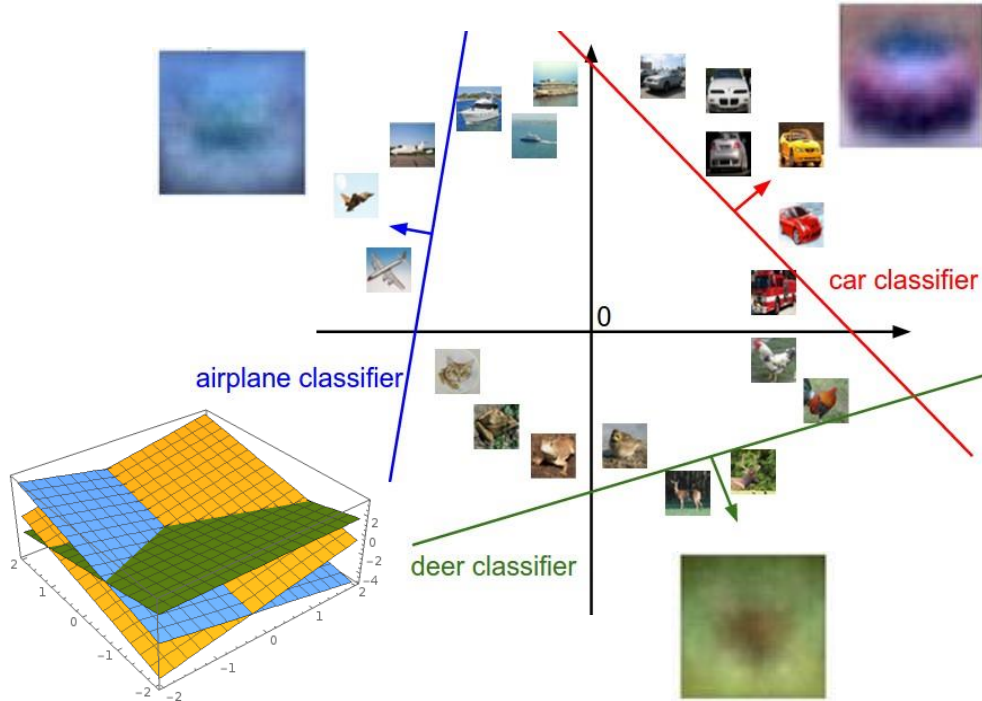
Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n

Geometric Viewpoint

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



Array of **32x32x3** numbers
(3072 numbers total)

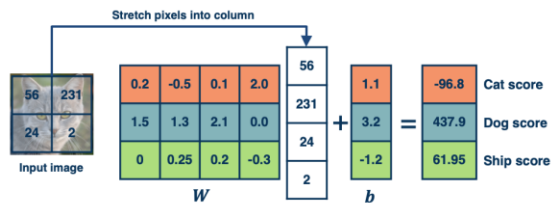


Plot created using Wolfram Cloud

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

Algebraic Viewpoint

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



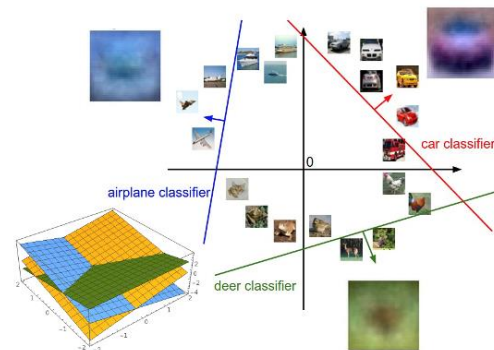
Visual Viewpoint

One template per class



Geometric Viewpoint

Hyperplanes cutting up space



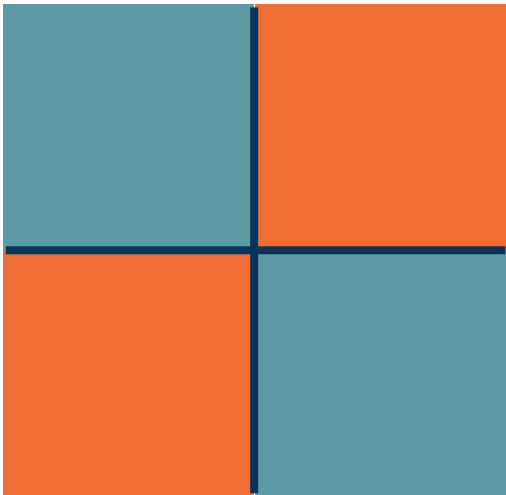
Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n

Class 1:

number of pixels > 0 odd

Class 2:

number of pixels > 0 even



Class 1:

$1 \leq \text{L2 norm} \leq 2$

Class 2:

Everything else

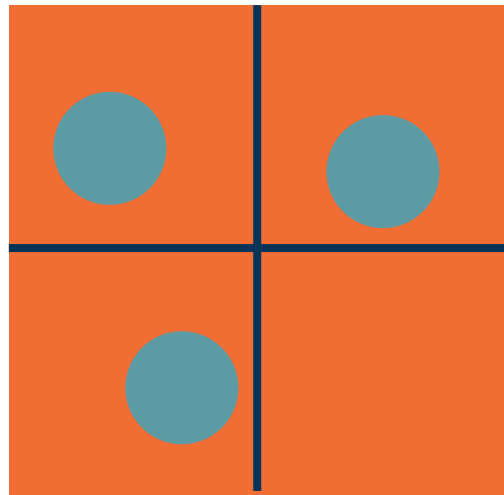


Class 1:

Three modes

Class 2:

Everything else



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

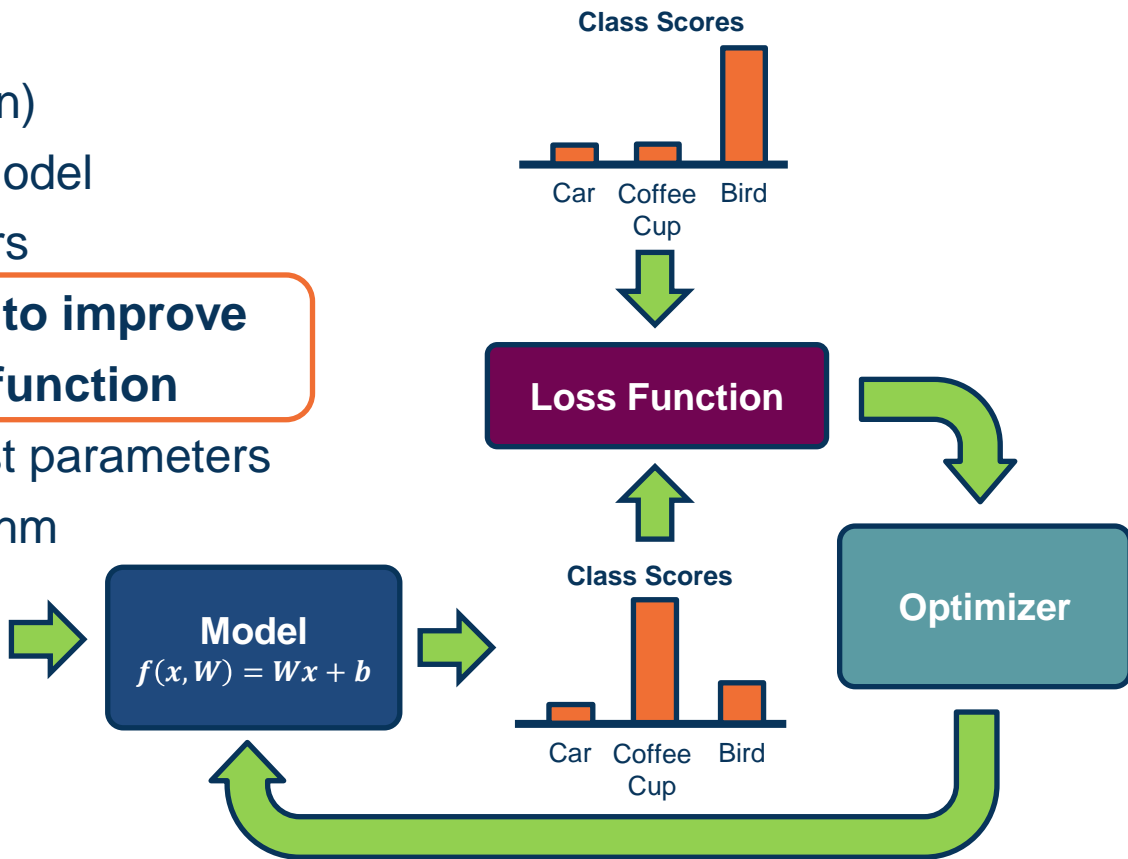
- We will learn complex, parameterized functions
 - Start w/ simple building blocks such as linear classifiers
- Key is to learn parameters, but learning is hard
 - Sources of generalization error
 - Add bias/assumptions via architecture, loss, optimizer
- Components of parametric classifiers:
 - Input/Output, Model (function), Loss function, Optimizer
 - Example: Image/Label, Linear Classifier, Hinge Loss, ?

Next Time:

- 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



Components of a Parametric Model

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,

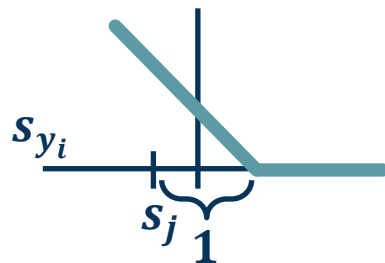


and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

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

Example: “Hinge Loss”



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

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:

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

$$\begin{aligned} &= \max(0, 5.1 - 3.2 + 1) \\ &\quad + \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 \\ &= 2.9 \end{aligned}$$

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

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:

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

$$\begin{aligned} &= \max(0, 1.3 - 4.9 + 1) \\ &\quad + \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

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:

0.0

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

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	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

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

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 is min/max of loss value?

[0, inf]



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

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 $s \approx 0$.
What is the loss?

C-1

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

Multiclass SVM loss:

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

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

No difference
(add constant 1)

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

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	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

- We will learn complex, parameterized functions
 - Start w/ simple building blocks such as linear classifiers
- Key is to learn parameters, but learning is hard
 - Sources of generalization error
 - Add bias/assumptions via architecture, loss, optimizer
- Components of parametric classifiers:
 - Input/Output, Model (function), Loss function, Optimizer
 - Example: Image/Label, Linear Classifier, Hinge Loss, ?

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

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

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

- 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

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

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

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

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

- 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)

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

Ground Truth

Prediction

$$\begin{aligned} \min_w KL(p^* || \hat{p}) &= \sum_y p^*(y) \log \frac{p^*(y)}{\hat{p}(y)} \\ &= \sum_y p^*(y) \log(p^*(y)) - \sum_y p^*(y) \log(\hat{p}(y)) \\ &\quad \underbrace{-H(p^*)}_{\text{(negative entropy, term goes away because not a function of model, } W, \text{ parameters we are minimizing over)}} \quad \underbrace{H(p^*, \hat{p})}_{\text{(Cross-Entropy)}} \end{aligned}$$

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

Softmax Classifier (Multinomial Logistic Regression)



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

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

Probabilities
must be ≥ 0

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

Softmax
Function

Probabilities
must sum to 1

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

cat
car
frog

3.2
5.1
-1.7

Unnormalized log-
probabilities / logits

exp

24.5
164.0
0.18

Unnormalized
probabilities

normalize

0.13
0.87
0.00

Probabilities



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

Q: How is it
possible that non-
GT probabilities
aren't in loss?

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**

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

Probabilities
must be ≥ 0

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

Softmax
Function

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**

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

Probabilities
must be ≥ 0

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

Softmax
Function

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),
 $-\log(1/C) = -\log(1) + \log(C)$
e.g. $\log(10) \approx 2$

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

Cross-Entropy Loss Example

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

L1 Regularization

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

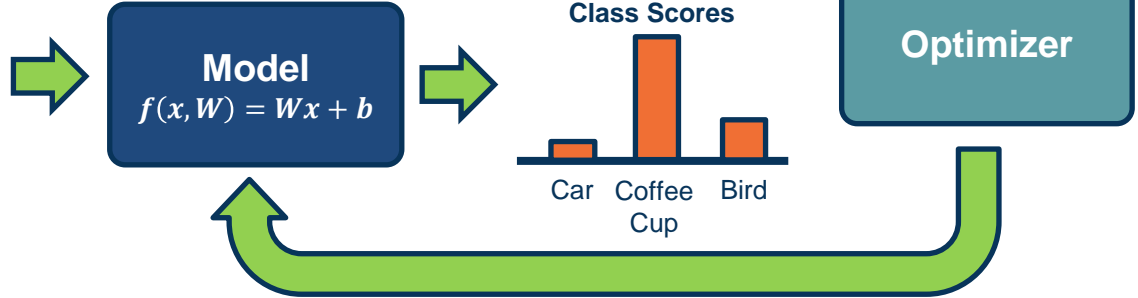
Example regularizations:

- L1/L2 on weights (encourage small values)

- 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



Components of a Parametric Model

Gradient Descent

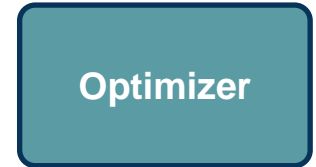
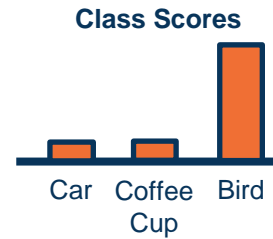
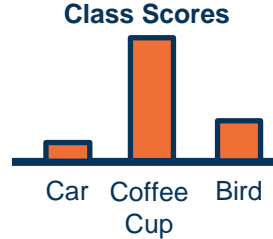
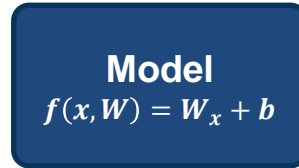
- Input (and representation)
- Functional form of the model
 - Including parameters
- Performance measure to improve
 - Loss or objective function
- **Algorithm for finding best parameters**
 - **Optimization algorithm**



Data: Image



Features: Histogram



Components of a Parametric Model

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

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b3 \end{bmatrix}$$

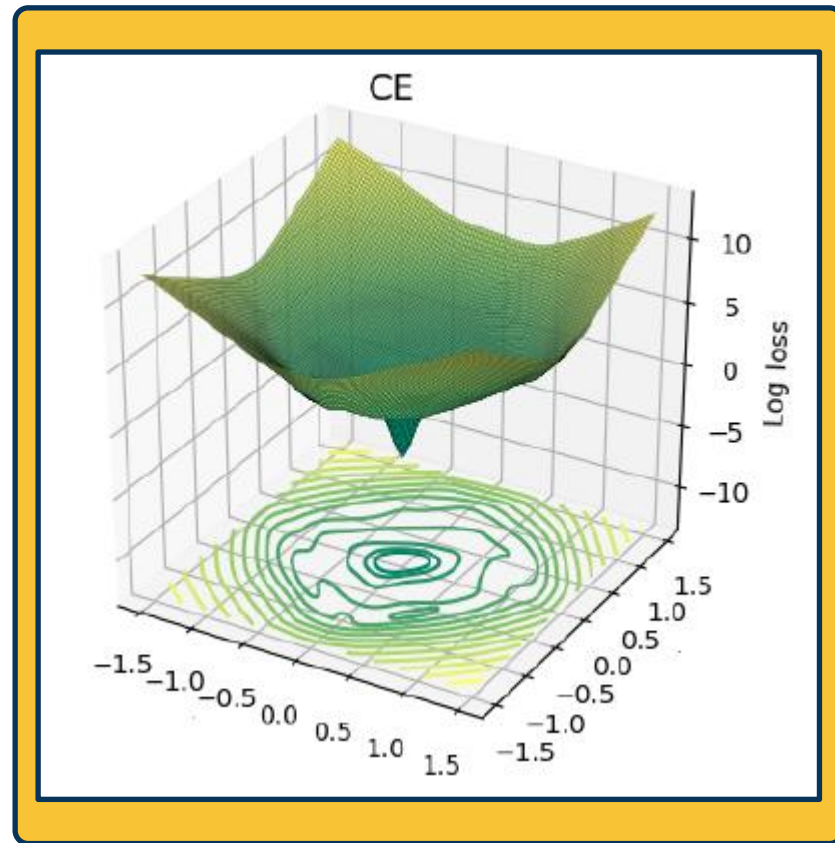


Loss

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**



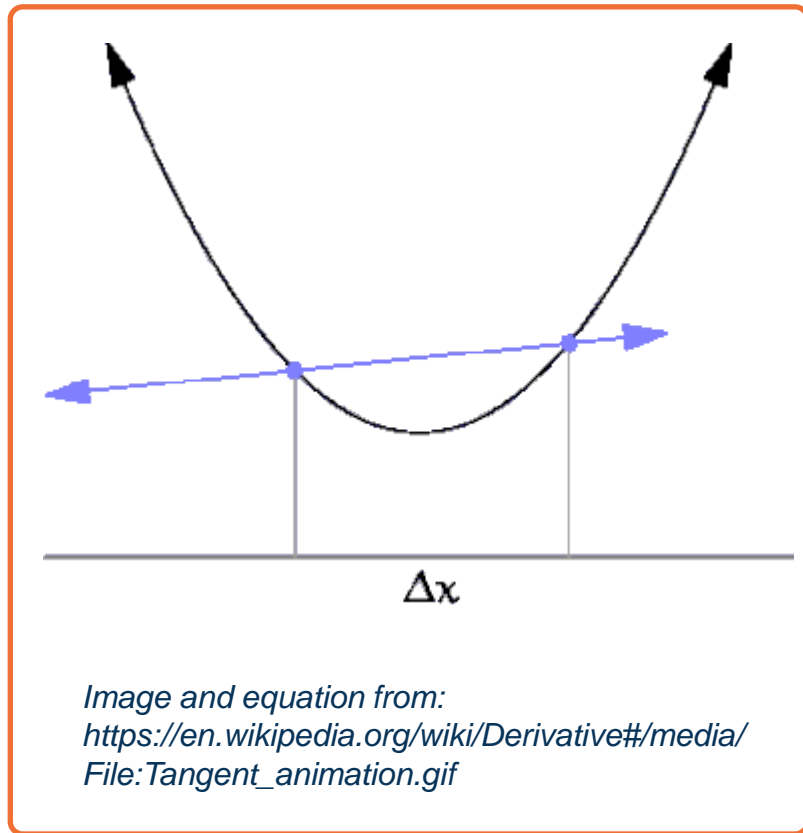


Strategy: Follow the Slope!

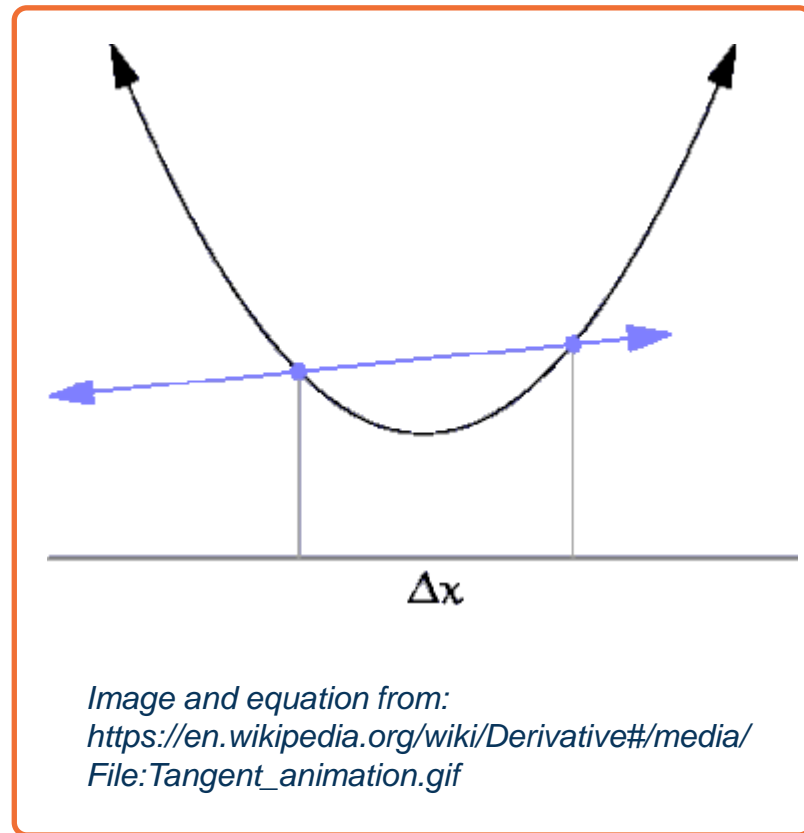
- We can find the steepest descent direction by computing the **derivative (gradient)**:

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

- 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



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

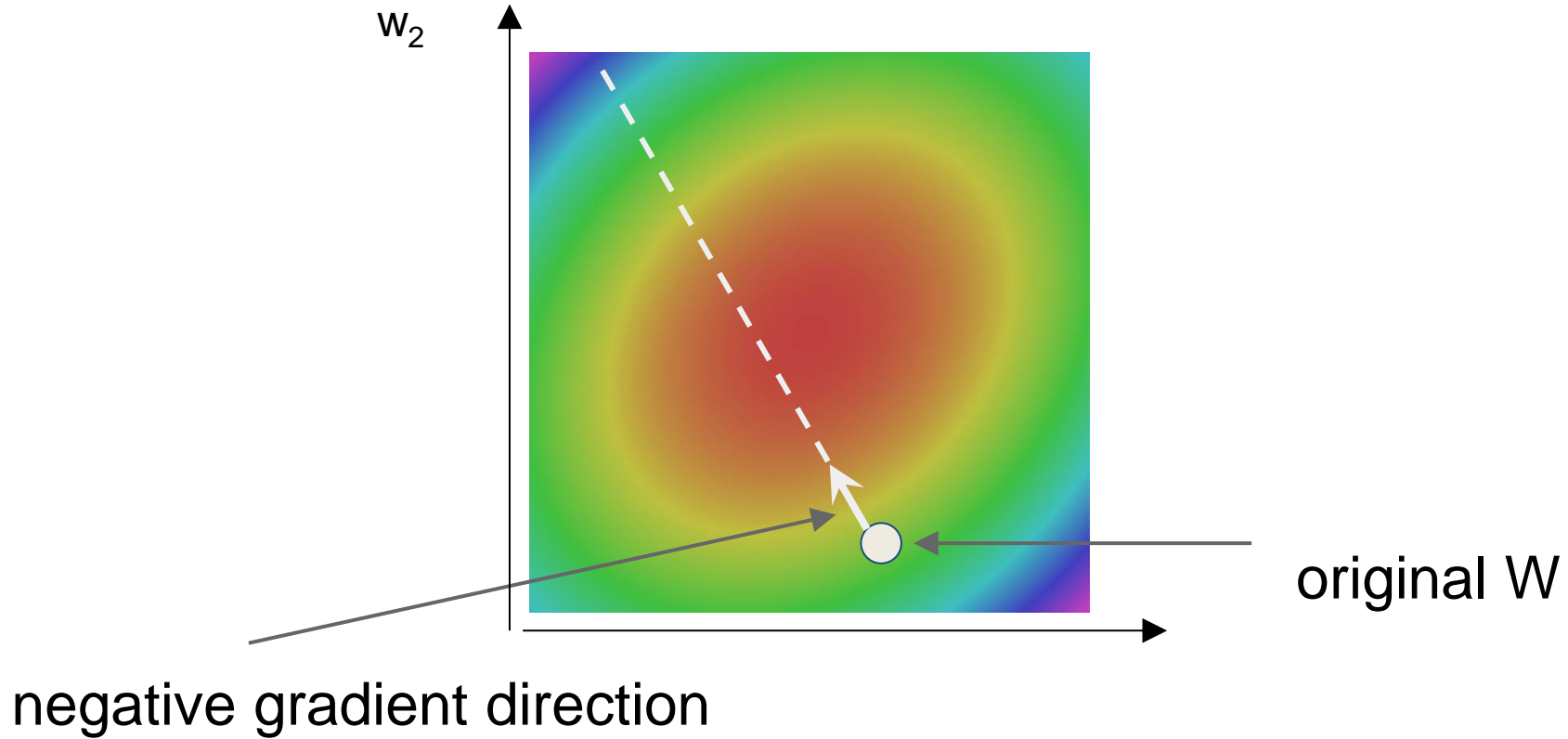


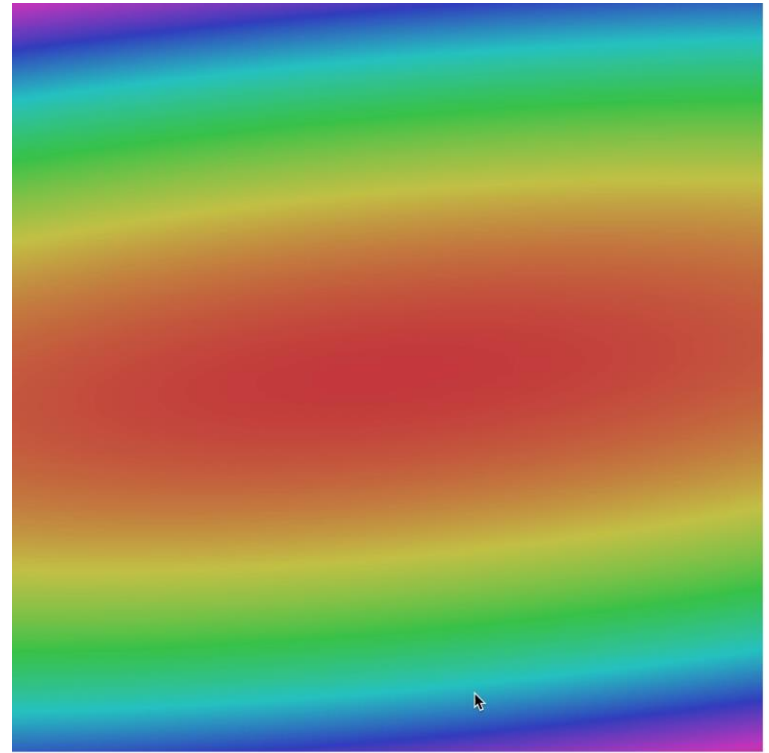
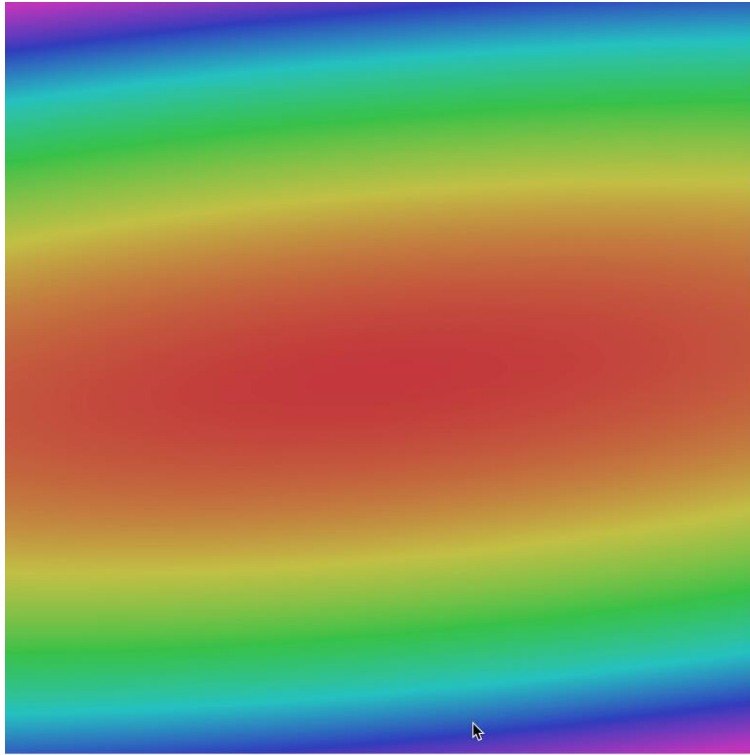
This idea can be turned into an **algorithm (gradient descent)**

1. Choose a model: $f(x, W) = Wx$
2. Choose loss function: $L_i = (y - Wx_i)^2$
3. Calculate partial derivative for each parameter: $\frac{\partial L}{\partial w_i}$
4. Update the parameters: $w_i = w_i - \frac{\partial L}{\partial w_i}$

Instead: Add learning rate to prevent too big of a step: $w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$

5. Repeat (from Step 3)





Gradient Descent

W_1

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

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

Mini-Batch Gradient Descent
$$L = \frac{1}{M} \sum L(f(x_i, W), y_i)$$

Where M is a *subset* of data

We iterate over mini-batches:

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

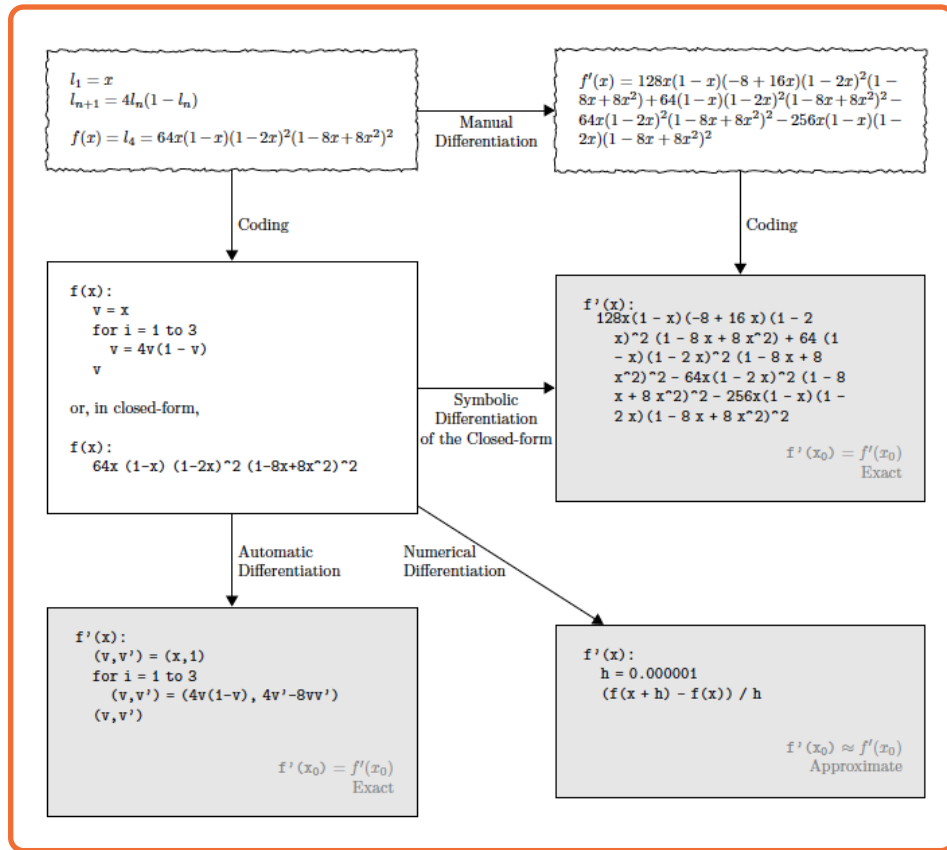
Gradient descent is guaranteed to converge under some conditions

- ✦ 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!

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



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,
?,
?,

$$(1.25322 - 1.25347)/0.0001 = -2.5$$

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

?,
?,...]

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,
?,
?,


$$(1.25353 - 1.25347) / 0.0001 = 0.6$$

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

?,...]

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,
?,
?,...]


$$(1.25347 - 1.25347)/0.0001 = 0$$

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

?,...]

Numerical vs Analytic Gradients

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

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

For some functions, we can analytically derive the partial derivative

Example:

Function

$$f(\mathbf{w}, \mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

(Assume \mathbf{w} and \mathbf{x}_i are column vectors, so same as $\mathbf{w} \cdot \mathbf{x}_i$)

Loss

$$\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Dataset: N examples (indexed by i)

Update Rule

$$\mathbf{w}_j \leftarrow \mathbf{w}_j + 2\alpha \sum_{i=1}^N \delta_i \mathbf{x}_{ij}$$

Derivation of Update Rule

$$L = \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Gradient descent tells us we should update \mathbf{w} as follows to minimize L :

$$\mathbf{w}_j \leftarrow \mathbf{w}_j - \alpha \frac{\partial L}{\partial \mathbf{w}_j}$$

So what's $\frac{\partial L}{\partial \mathbf{w}_j}$?

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}_j} &= \sum_{i=1}^N \frac{\partial}{\partial \mathbf{w}_j} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \\ &= \sum_{i=1}^N 2(y_i - \mathbf{w}^T \mathbf{x}_i) \frac{\partial}{\partial \mathbf{w}_j} (y_i - \mathbf{w}^T \mathbf{x}_i) \\ &= -2 \sum_{i=1}^N \delta_i \frac{\partial}{\partial \mathbf{w}_j} \mathbf{w}^T \mathbf{x}_i \\ &\quad \dots \text{where} \dots \\ &\quad \delta_i = y_i - \mathbf{w}^T \mathbf{x}_i \\ &= -2 \sum_{i=1}^N \delta_i \frac{\partial}{\partial \mathbf{w}_j} \sum_{k=1}^N \mathbf{w}_k \mathbf{x}_{ik} \\ &= -2 \sum_{i=1}^N \delta_i \mathbf{x}_{ij} \end{aligned}$$

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

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

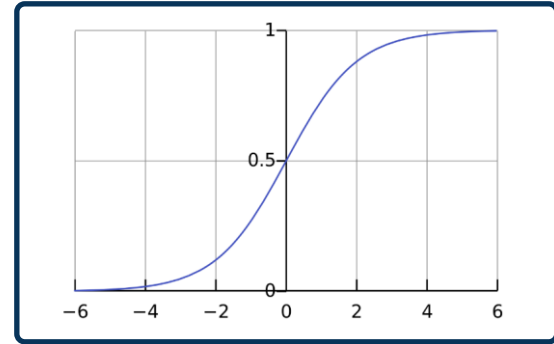
First, one can derive that: $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

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

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

$$\begin{aligned} \frac{\partial L}{\partial w_j} &= \sum_i 2 \left(y_i - \sigma\left(\sum_k w_k x_{ik}\right) \right) \left(-\frac{\partial}{\partial w_j} \sigma\left(\sum_k w_k x_{ik}\right) \right) \\ &= \sum_i -2 \left(y_i - \sigma\left(\sum_k w_k x_{ik}\right) \right) \sigma'\left(\sum_k w_k x_{ik}\right) \frac{\partial}{\partial w_j} \sum_k w_k x_{ik} \\ &= \sum_i -2 \delta_i \sigma(d_i) (1 - \sigma(d_i)) x_{ij} \end{aligned}$$

$$\text{where } \delta_i = y_i - \mathbf{f}(x_i) \quad d_i = \sum_k w_k x_{ik}$$



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

$$\delta_i = y_i - \sigma_i$$

- We will learn complex, parameterized functions
 - Start w/ simple building blocks such as linear classifiers
- Optimize parameters via simple gradient descent (!)
- But calculating the gradients is cumbersome for more complex functions
- Let's develop a generic representation of these functions and an algorithm that can do this easily!