

# Introduction to Machine Learning (Basic Concepts)

CS5330, Huaizu Jiang

Fall 2021, Northeastern University

# Before the class

- GPU resources for the final project
  - Discovery (**you should have received notifications**)
  - AWS
  - Google could
- Difficulty of pa3 vs pa2

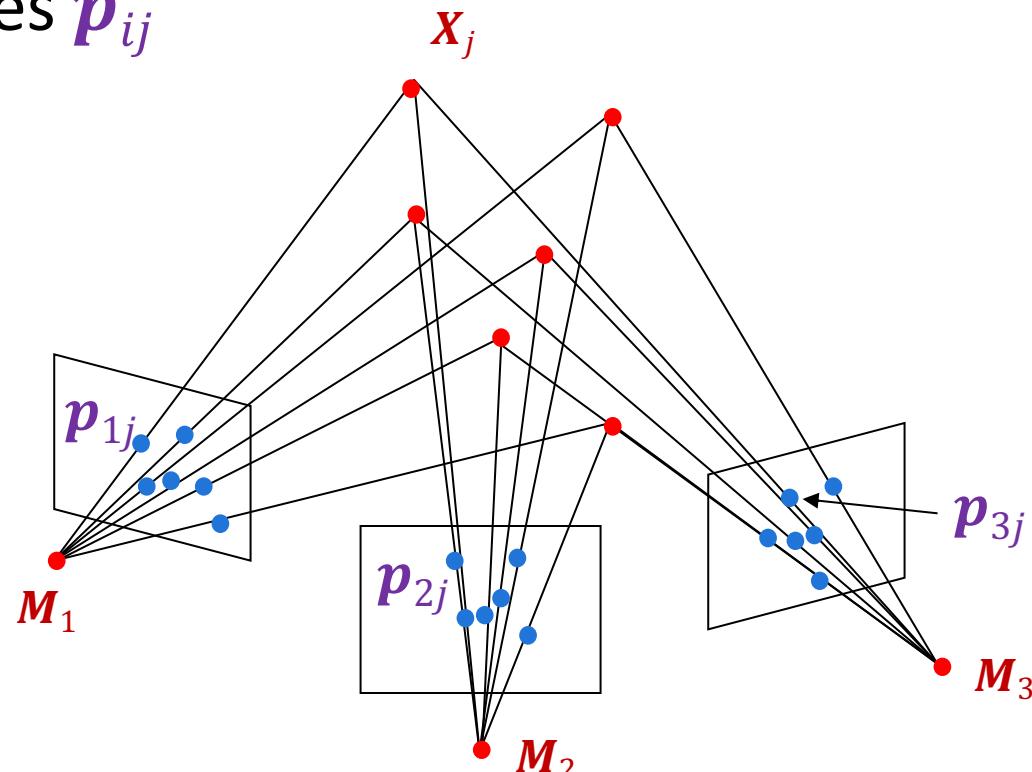
# Recap

# Structure from motion: Problem formulation

- Given:  $m$  images of  $n$  fixed 3D points such that (ignoring visibility)

- $\bullet \ p_{ij} \equiv \mathbf{M}_i \mathbf{X}_j, \quad i = 1, \dots, m, \ j = 1, \dots, n$

- Problem:** estimate  $m$  projection matrices  $\mathbf{M}_i$  and  $n$  3D points  $\mathbf{X}_j$  from the  $mn$  correspondences  $p_{ij}$



# Structure from motion ambiguities

Not just limited to scale. Given:

$$\mathbf{p}_{ij} \equiv \mathbf{M}_i \mathbf{X}_j$$

Can insert any global transform  $\mathbf{H}$

$$\mathbf{p}_{ij} \equiv \mathbf{M}_i \mathbf{X}_j = \mathbf{M}_i \mathbf{H}^{-1} \mathbf{H} \mathbf{X}_j$$

$\mathbf{H}$  is a 3D homography / perspective transform / projective transform

# Affine structure from motion

General structure  
from motion:

$$\mathbf{p}_{ij} \equiv M_i \mathbf{X}_j$$

3x1            3x4    4x1

Assume M is affine  
camera:

$$\mathbf{p}_{ij} = A_i \mathbf{X}_j + \mathbf{b}_i$$

2x1            2x3    3x1            2x1

- **Problem:** use the  $mn$  2D points  $\mathbf{p}_{ij}$  to estimate  $m$  projection matrices  $A_i$  and translation vectors  $b_i$ , and  $n$  3D points  $X_j$

# Using Normalized 2D Points

- After centering, each normalized 2D point  $\hat{p}_{ij}$  is related to the normalized 3D point by

$$\bullet \hat{p}_{ij} = A_i \hat{X}_j$$

- We can get rid of the need to center the 3D data (and the translation ambiguity) by defining the origin of the world coordinate system as the centroid of the 3D points

# Affine structure from motion

Then, write all the equations in one in terms of product of cameras and points.

$$\widehat{\mathbf{p}_{ij}} = \mathbf{A}_i \widehat{\mathbf{X}}_j \quad \mathbf{D} = \begin{bmatrix} \widehat{\mathbf{p}_{11}} & \cdots & \widehat{\mathbf{p}_{1n}} \\ \vdots & \ddots & \vdots \\ \widehat{\mathbf{p}_{m1}} & \cdots & \widehat{\mathbf{p}_{mn}} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_m \end{bmatrix} [\mathbf{X}_1 \quad \cdots \quad \mathbf{X}_n]$$

2m x n

D

2mx3

M

3xn

S

What's the rank of D?

3!

C. Tomasi and T. Kanade. [Shape and motion from image streams under orthography: A factorization method.](#) IJCV, 9(2):137-154, November 1992.

# Factorizing the measurement matrix

- Keep top 3 singular values:
  - This is the closest approximation of  $\mathbf{D}$  with a rank-3 matrix in terms of Frobenius norm

$$\mathbf{D}_{2m \times n} = \mathbf{U}_3_{2m \times 3} \times \begin{matrix} \Sigma_3 \\ 3 \times 3 \end{matrix} \times \mathbf{V}_3^T_{n \times n}$$

- What to do about  $\Sigma_3$ ?
- One solution:  $\mathbf{M} = \mathbf{U}_3 \Sigma_3^{\frac{1}{2}}, \mathbf{S} = \Sigma_3^{\frac{1}{2}} \mathbf{V}_3^T$

# Factorizing the measurement matrix

- One possible solution:

$$\begin{matrix} D \\ 2m \times n \end{matrix} = \begin{matrix} M \\ 2m \times 3 \end{matrix} \times \begin{matrix} S \\ 3 \times n \end{matrix}$$

$$S = \Sigma_3^{\frac{1}{2}} V_3^T$$

$$M = U_3 \Sigma_3^{\frac{1}{2}}$$

- Are there other solutions?

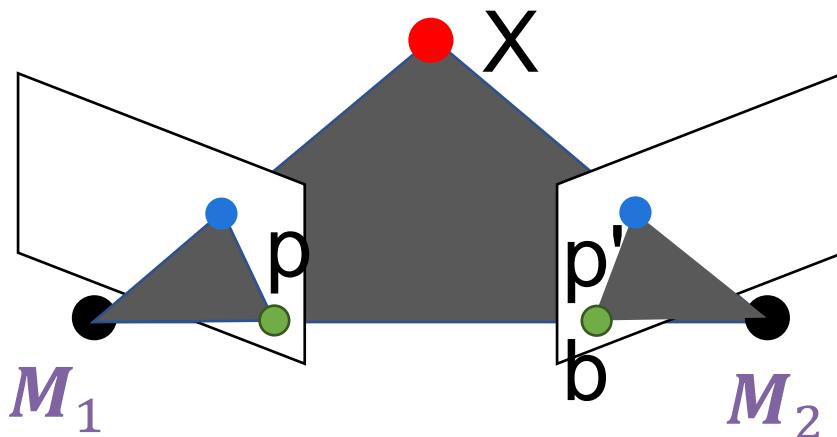
# Projective structure from motion

- Given:  $m$  images of  $n$  fixed 3D points such that (ignoring visibility):
  - $\mathbf{p}_{ij} \equiv \mathbf{M}_i \mathbf{X}_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$
- Problem: estimate  $m$  projection matrices  $\mathbf{M}_i$  and  $n$  3D points  $\mathbf{X}_j$  from the  $mn$  correspondences  $\mathbf{p}_{ij}$
- With no calibration info, cameras and points can only be recovered up to a  $4 \times 4$  projective transformation  $\mathbf{Q}$ :
  - $\mathbf{X} \rightarrow \mathbf{Q}\mathbf{X}, \quad \mathbf{M} \rightarrow \mathbf{M}\mathbf{Q}^{-1}$

# Two Camera Case

For two cameras, we need 7 points.

**What else (in theory) requires 7 points?**



Compute fundamental matrix  $\mathbf{F}$  and epipole  $\mathbf{b}$  s.t.  $\mathbf{F}^T \mathbf{b} = 0$ . Then:

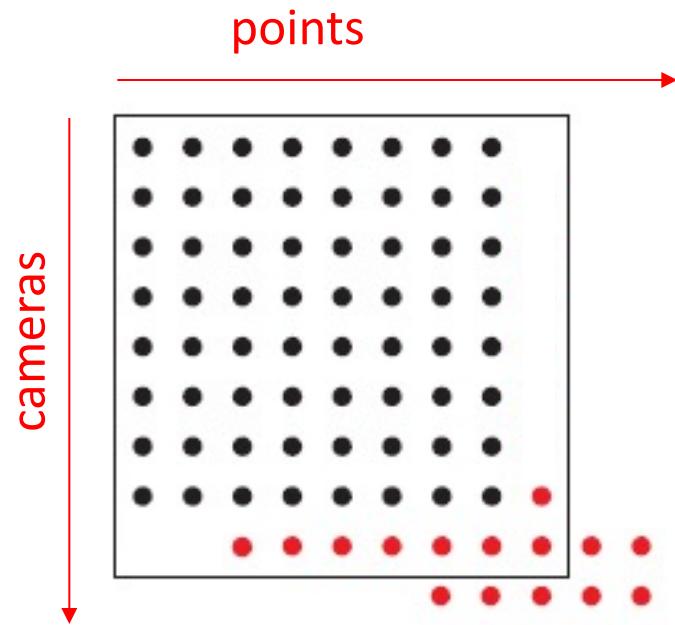
$$\mathbf{M}_1 = [\mathbf{I}, \mathbf{0}]$$

$$\mathbf{M}_2 = [-[\mathbf{b}_x]\mathbf{F}, \mathbf{b}]$$

Remember: this is up to a projective ambiguity!

# Incremental structure from motion

- Initialize motion from two images using fundamental matrix
- Initialize structure by triangulation
- For each additional view:
  - Determine projection matrix of new camera using all the known 3D points that are visible in its image – **calibration**
  - Refine and extend structure: compute new 3D points, re-optimize existing points that are also seen by this camera – **triangulation**
- Refine structure and motion: bundle adjustment

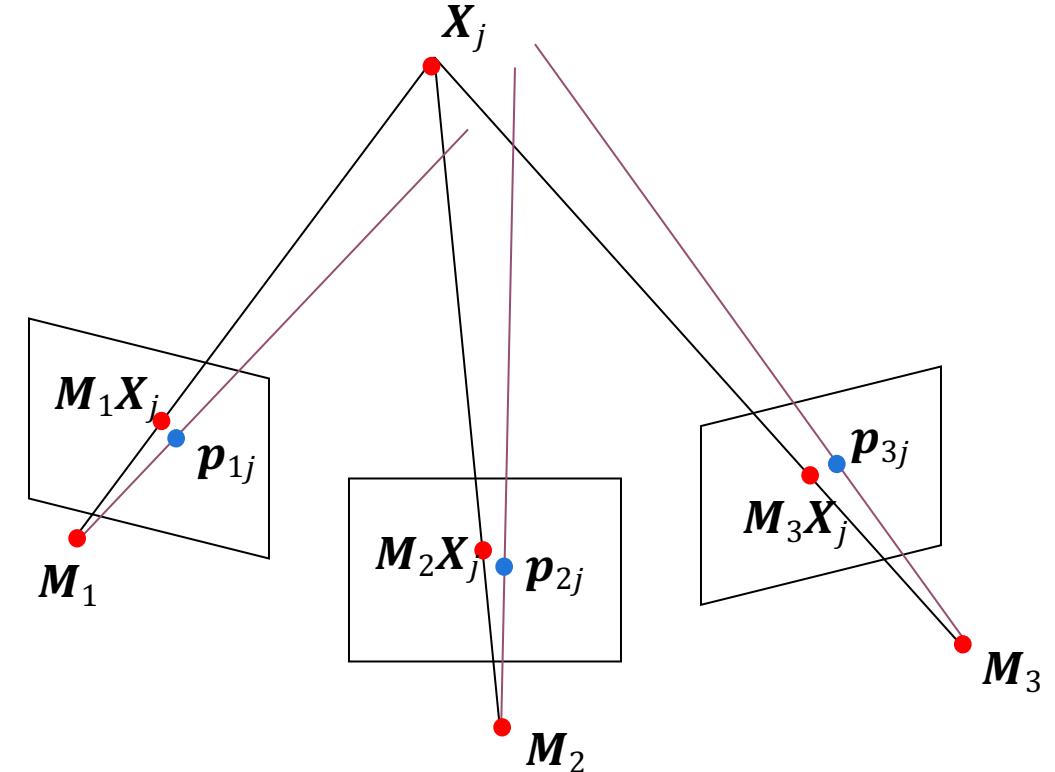


# Bundle adjustment

- Non-linear method for refining structure and motion
- Minimize reprojection error (with lots of bells and whistles):

$$\bullet \sum_{i=1}^m \sum_{j=1}^n w_{ij} d(\mathbf{p}_{ij} - \mathbf{M}_i \mathbf{X}_j)^2$$

↑  
visibility flag: is  
point  $j$  visible in  
view  $i$ ?



B. Triggs et al. [Bundle adjustment – A modern synthesis](#). International Workshop on Vision Algorithms, 1999

# Today's class

- A bit more about Structure from motion (SfM)
- Introduction to machine learning
  - What is machine learning
  - A simple machine learning example
  - Cross validation

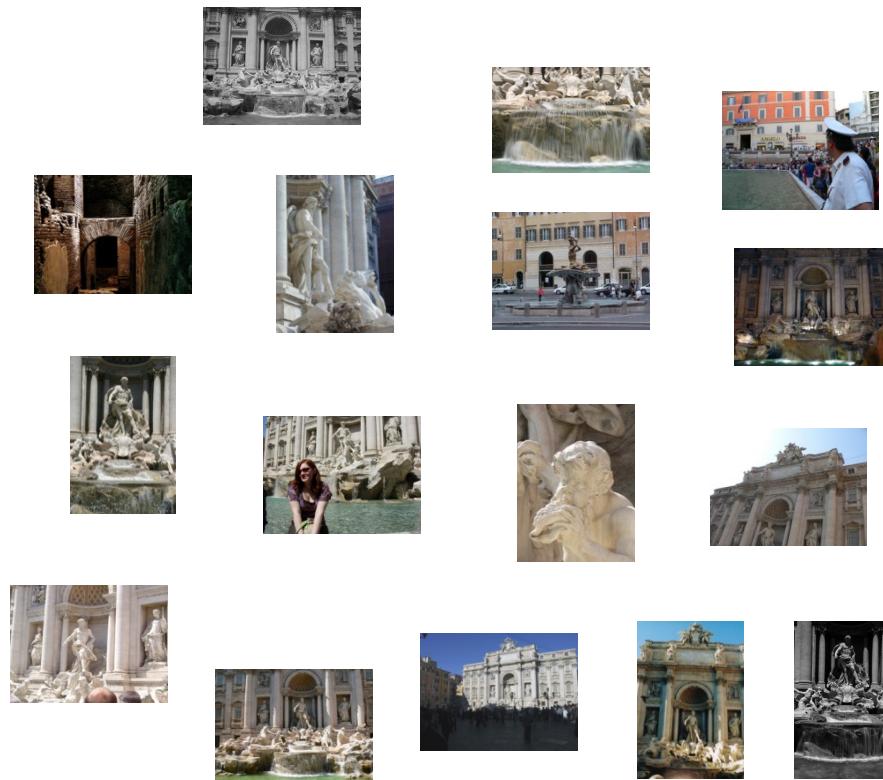
# Representative SFM pipeline



N. Snavely, S. Seitz, and R. Szeliski. [Photo tourism: Exploring photo collections in 3D](#). SIGGRAPH 2006  
<http://phototour.cs.washington.edu/>

# Feature detection

## Detect SIFT features



Source: N. Snavely

# Feature detection

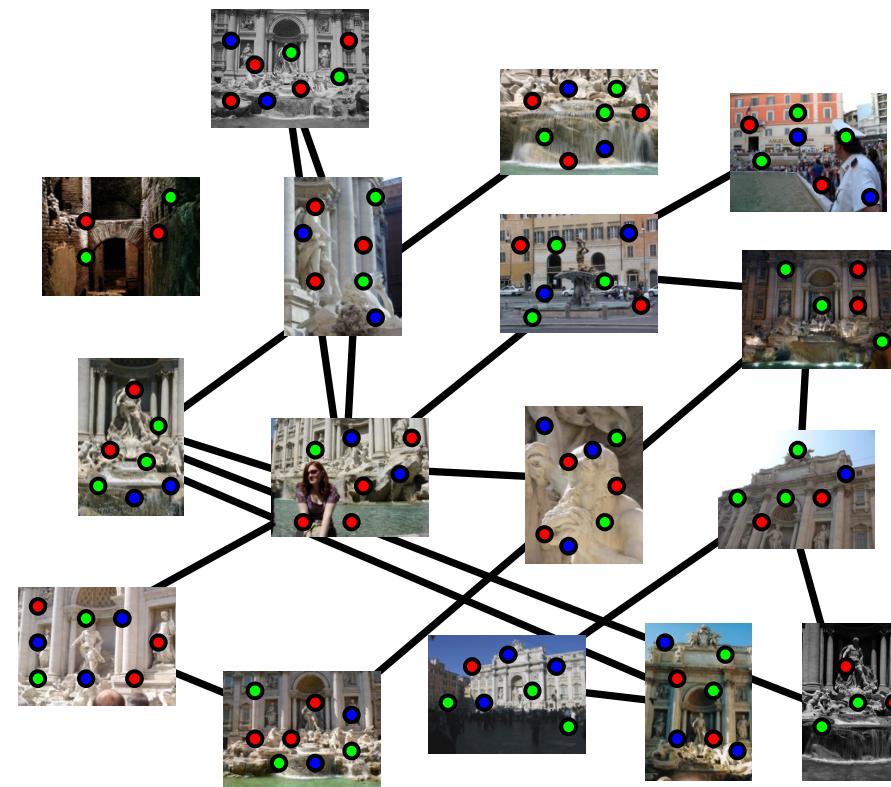
## Detect SIFT features

Other popular feature types: [SURF](#), [ORB](#), [BRISK](#), ...



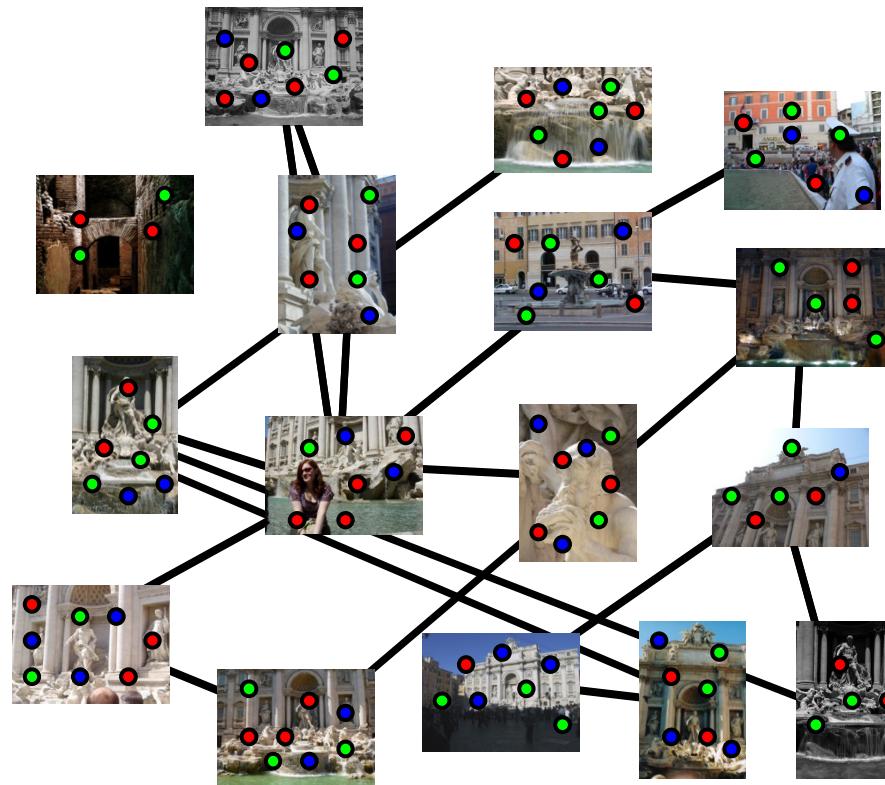
# Feature matching

Match features between each pair of images

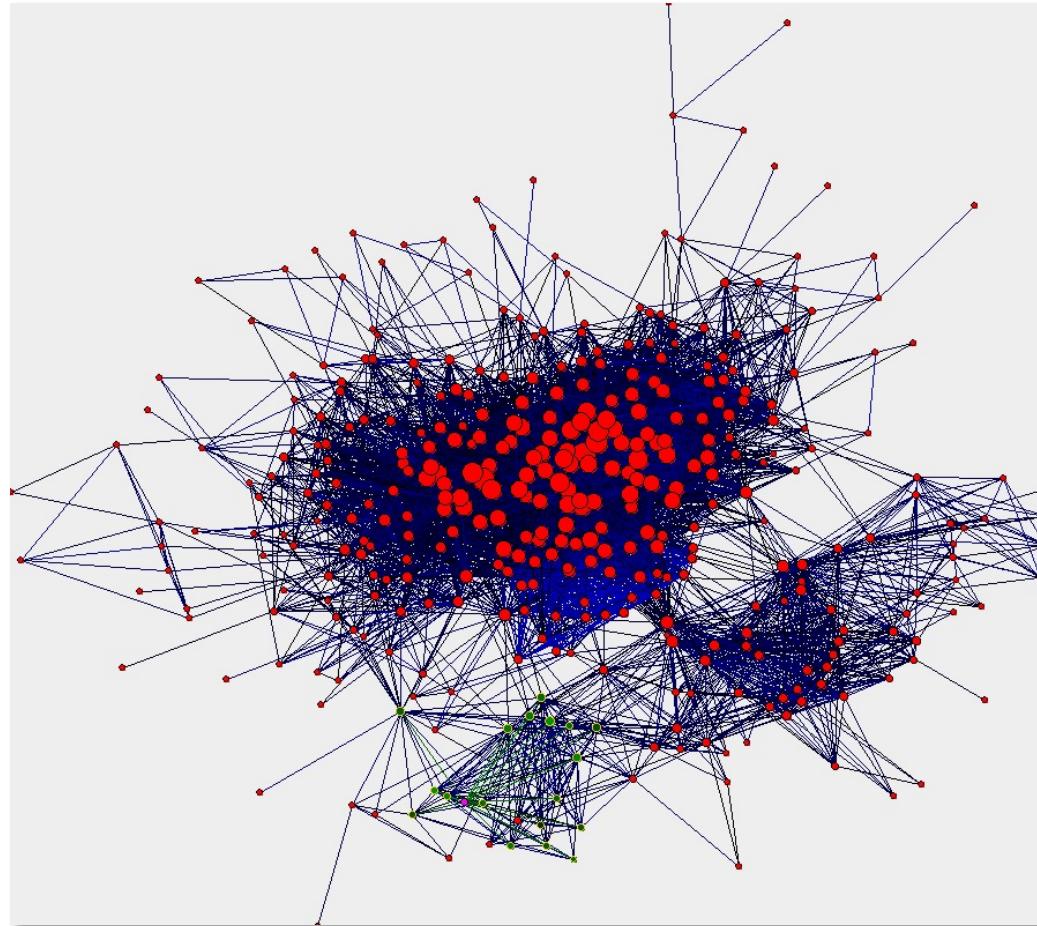


# Feature matching

Use RANSAC to estimate fundamental matrix between each pair



# Image connectivity graph



(graph layout produced using the Graphviz toolkit: <http://www.graphviz.org/>)

# Incremental SFM

- Pick a pair of images with lots of inliers (and preferably, good EXIF data)
  - Initialize intrinsic parameters (focal length, principal point) from EXIF
  - Estimate extrinsic parameters ( $R$  and  $t$ ) using [five-point algorithm](#)
  - Use triangulation to initialize model points
- While remaining images exist
  - Find an image with many feature matches with images in the model
  - Run RANSAC on feature matches to register new image to model
  - Triangulate new points
  - Perform bundle adjustment to re-optimize everything
  - Optionally, align with GPS from EXIF data or ground control points

# SfM applications

- 3D modeling
- Surveying
- Robot navigation and mapmaking (**including autonomous driving, SLAM**)
- Virtual and augmented reality
- Visual effects ("Match moving")
  - [https://www.youtube.com/watch?v=RdYWp70P\\_kY](https://www.youtube.com/watch?v=RdYWp70P_kY)

# Applications: Visual Reality & Augmented Reality & Metaverse



Oculus

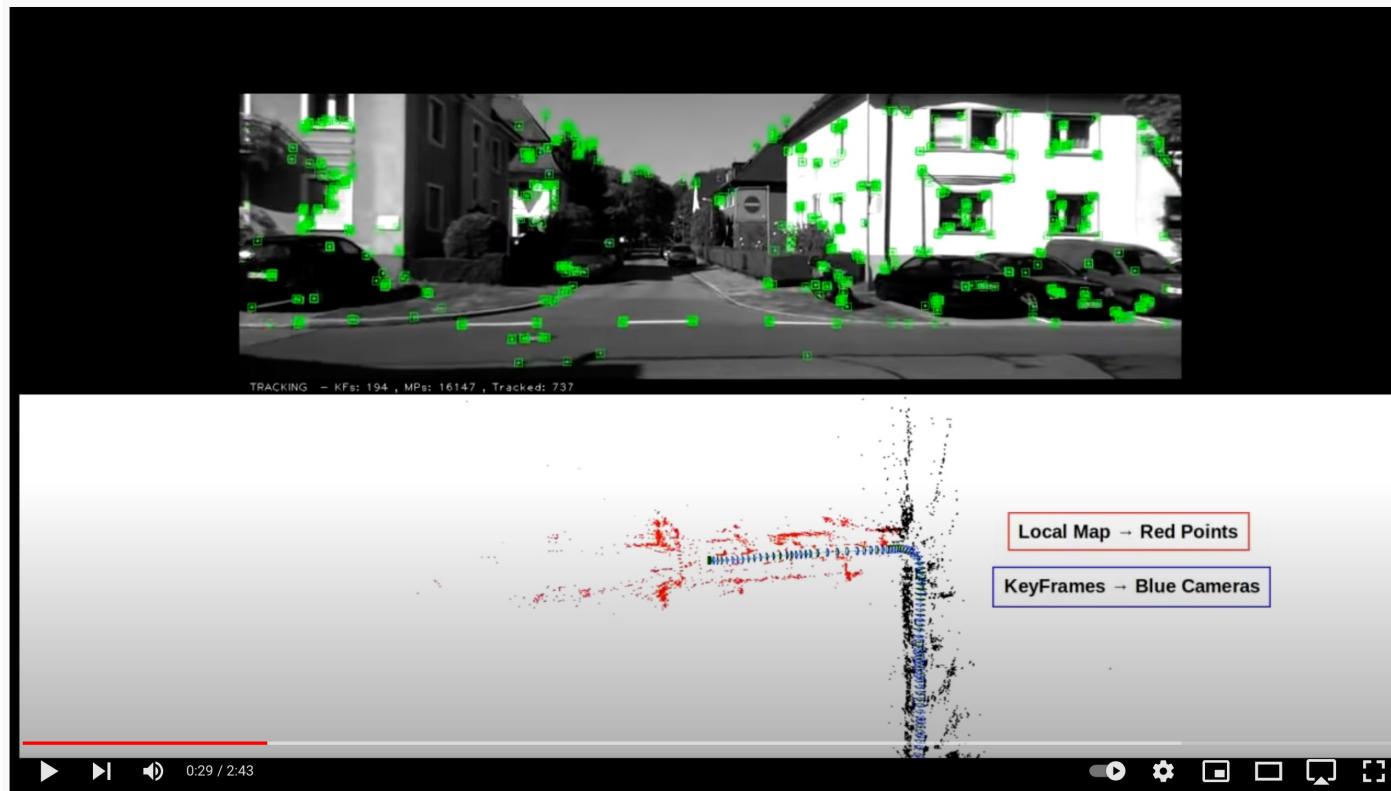
<https://www.youtube.com/watch?v=KOG7yTz1iTA>



HoloLens

<https://www.youtube.com/watch?v=FMtvrTGnP04>

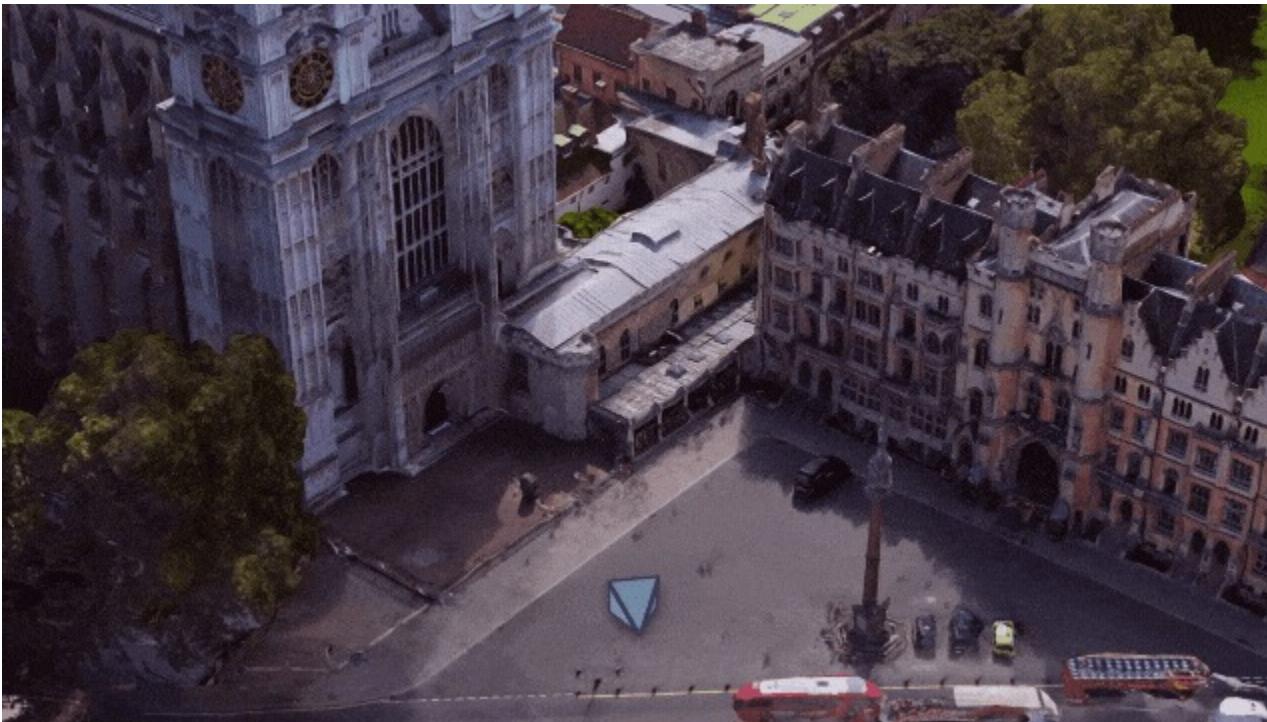
# Applications: Robot Navigation (Autonomous Driving)



ORB-SLAM

<https://www.youtube.com/watch?v=8DISRmsO2YQ>

# Application: AR walking directions



<https://www.theverge.com/2019/8/8/20776247/google-maps-live-view-ar-walking-directions-ios-android-feature>

# The rest of this course

<b>Fundamentals of Deep Neural Networks</b>		
Week 7	Mon Oct 25	Basic Concepts of Machine Learning
	Thu, Oct 28	Optimization: Stochastic Gradient Descent and Backpropagation
Week 8	Mon, Nov 1	Neural Networks 1
	Thu, Nov 4	Neural Networks 2
Week 9	Mon, Nov 8	Convolutional Neural Networks 1
	Thu, Nov 11	No class, Veterans' Day
Week 10	Mon, Nov 15	Convolutional Neural Networks 2
	<b>Recognition, Detection and Segmentation via DNNs</b>	
Week 11	Thur, Nov 18	Object Detection 1: R-CNN
	Mon, Nov 22	Object Detection 2: Fast and Faster R-CNN
Week 12	Thur, Nov 25	No class, Thanksgiving recess
	Mon, No 29	Segmentation 1: Mask R-CNN, FCN
Week 13	Thur, Dec 2	Segmentation 2: Panoptic Segmentation
	<b>Advanced Topics</b>	
Week 14	Mon, Dec 6	Novel view synthesis and video frame interpolation
	Thur, Dec 9	Vision and language (VQA, visual grounding, etc)
	Mon, Dec 13	Learning with unlabeled data
	Thur, Dec 16	Transformers in computer vision

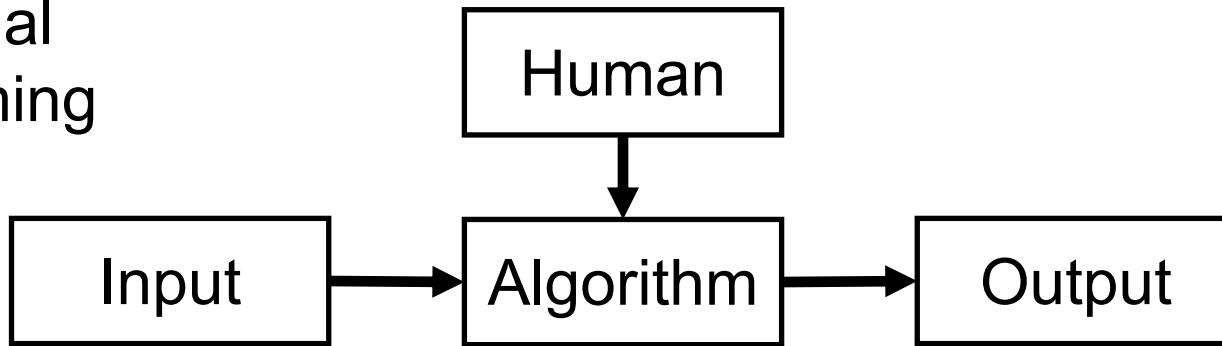
guest lecture

# Machine Learning

Algorithms that learn from data

# Machine Learning vs Programming

Traditional  
Programming

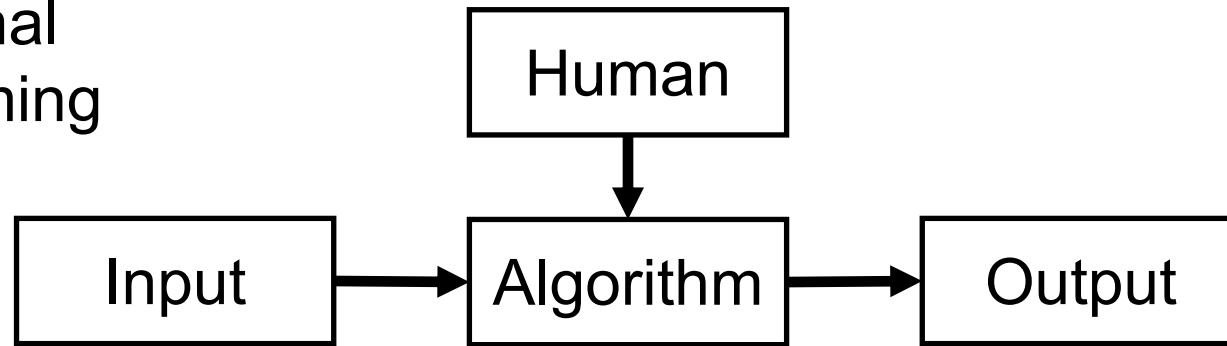


Works well for  
sorting  
numbers

```
def bubble_sort(arr):
    N = len(arr)
    for i in range(N - 1):
        for j in range(N - i - 1):
            if arr[j + 1] < arr[j]:
                temp = arr[j]
                arr[j] = arr[j + 1]
                arr[j + 1] = temp
    return arr
```

# Machine Learning vs Programming

Traditional  
Programming



Much harder  
for some  
problems

```
def cat_or_dog(image):  
    if ???:  
        return "cat"  
    else:  
        return "dog"
```

# Image Classification: a core task in Computer Vision



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



cat

# The problem: *semantic gap*

Images are represented as  
3D arrays of numbers, with  
integers between [0, 255].

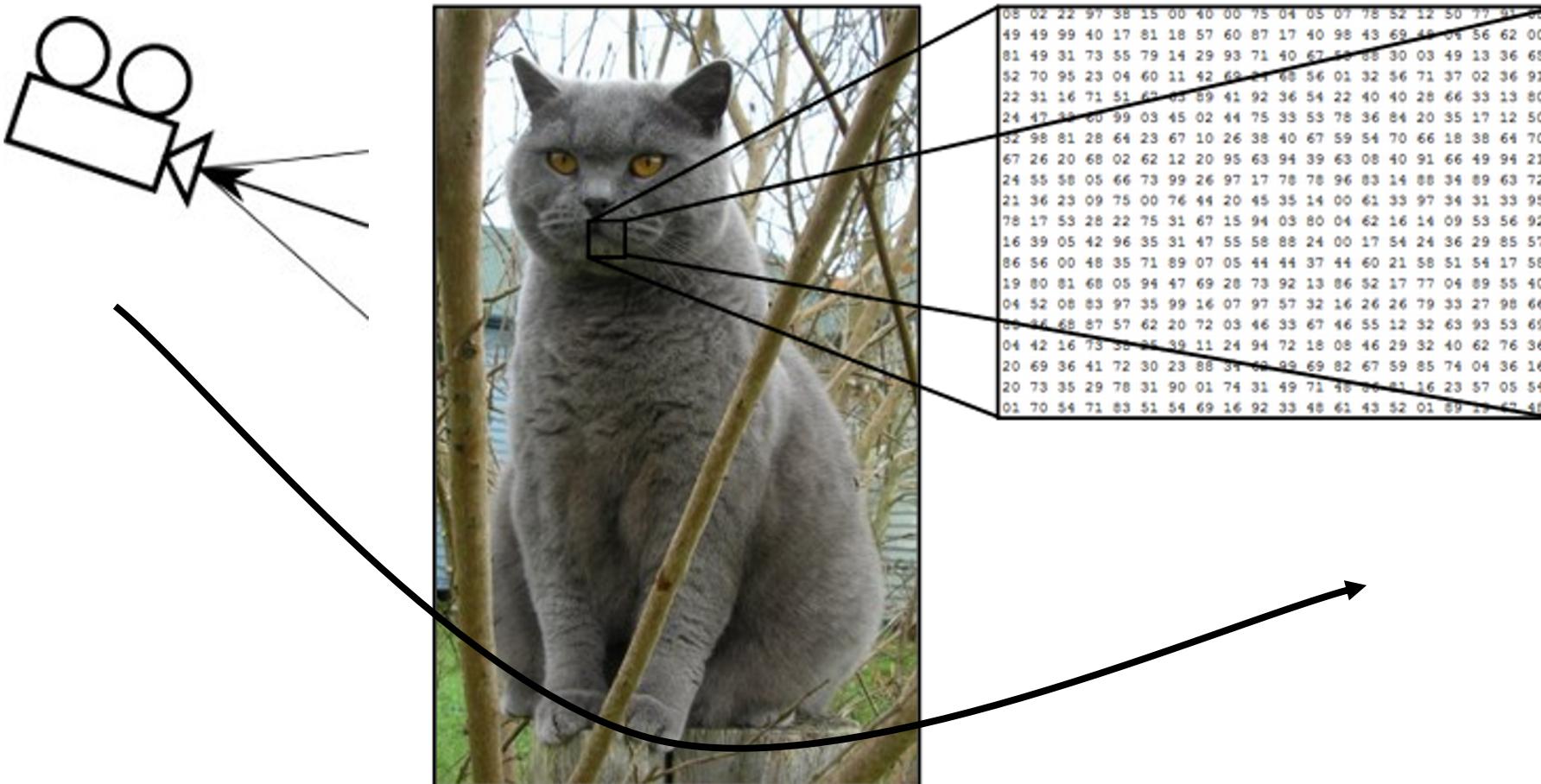
E.g.  
300 x 100 x 3



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	91	68
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	53	68	30	03	49	13	36	65
52	70	95	23	04	60	11	42	62	21	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	63	59	41	92	36	54	22	40	40	28	66	33	13	80
24	47	81	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
05	44	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	30	25	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	34	02	99	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	48	86	81	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	17	67	48

What the computer sees

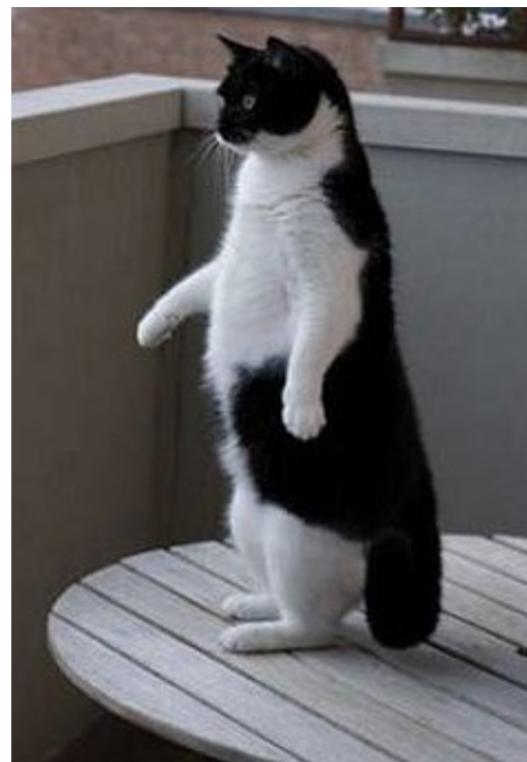
# Challenges: Viewpoint Variation



# Challenges: Illumination



# Challenges: Deformation



# Challenges: Occlusion



# Challenges: Background clutter

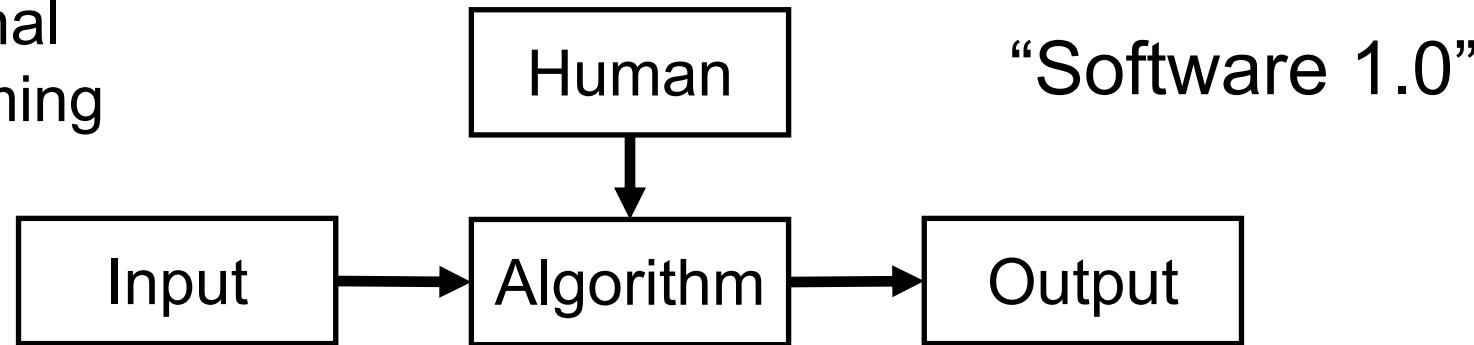


# Challenges: Intraclass variation

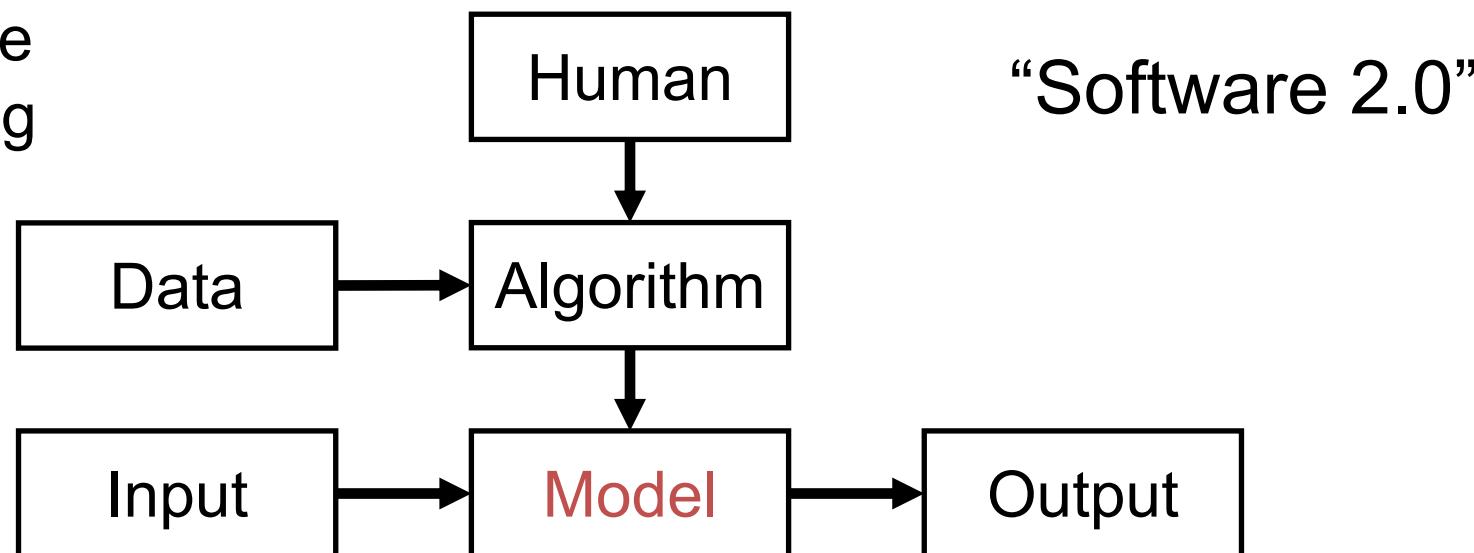


# Machine Learning vs Programming

Traditional  
Programming



Machine  
Learning



<https://medium.com/@karpathy/software-2-0-a64152b37c35>

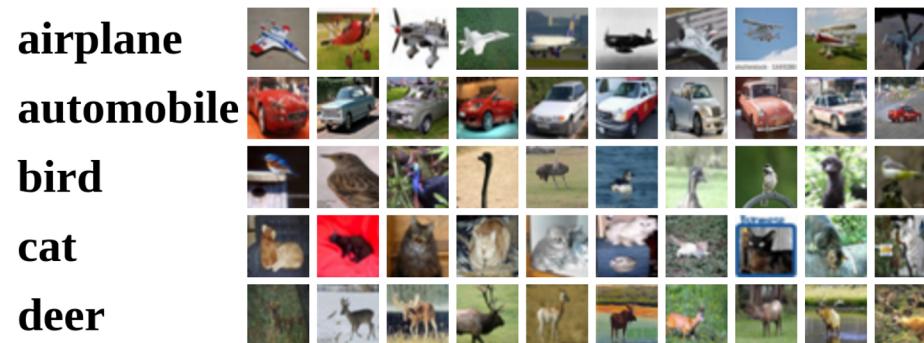
# Machine Learning: Data-Driven Approach

1. Collect a large set of data
2. Use Machine Learning to train a model
3. Evaluate the model on new data

```
def train(images, labels):  
    # Machine learning!  
    return model
```

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

**Example training set**



# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is input / feature

$y$  is label / target

**Goal:** Learn a  
*function* to map  $x \rightarrow y$

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is input / feature

$y$  is label / target

**Goal:** Learn a  
*function* to map  $x \rightarrow y$

Image Classification:  
Predict a discrete category

$x$



$y$

→ Cat



→ Dog



→ Monkey

Cat image is CC0 public domain  
Dog image is CC0 Public Domain  
Monkey image is CC0 Public Domain

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is input / feature

$y$  is label / target

**Goal:** Learn a  
*function* to map  $x \rightarrow y$

Image Regression:  
Predict a continuous value

$x$

$y$



→ 2 lbs



→ 25 lbs



→ 35 lbs

Cat image is CC0 public domain  
Dog image is CC0 Public Domain  
Monkey image is CC0 Public Domain

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is input / feature

$y$  is label / target

**Goal:** Learn a  
*function* to map  $x \rightarrow y$

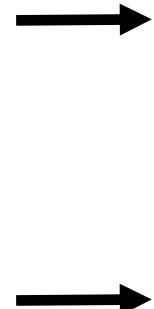
**Image Captioning:**  
Predict a sequence of words

$x$



$y$

“A white and gray kitten on grass”



“White and orange dog with a red leash in the woods”



“A monkey sitting in front of rocks”

Cat image is CC0 public domain

Dog image is CC0 Public Domain

Monkey image is CC0 Public Domain

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is input / feature

$y$  is label / target

**Goal:** Learn a  
*function* to map  $x \rightarrow y$

## Unsupervised Learning

**Data:**  $x$

Just data, no labels!

**Goal:** Learn  
underlying *structure*  
in the data

# Supervised vs Unsupervised Learning

Clustering:  
Group similar images



Cat image is CC0 public domain  
Dog image is CC0 Public Domain  
Monkey image is CC0 Public Domain

Cat image is CC0 public domain  
Dog image is CC0 public domain  
Monkey image is CC0 public domain

**Unsupervised  
Learning**

**Data:** x  
Just data, no labels!

**Goal:** Learn  
underlying *structure*  
in the data

# Supervised vs Unsupervised Learning

Dimensionality Reduction:  
Project to subspace

Images:  
256x256x3

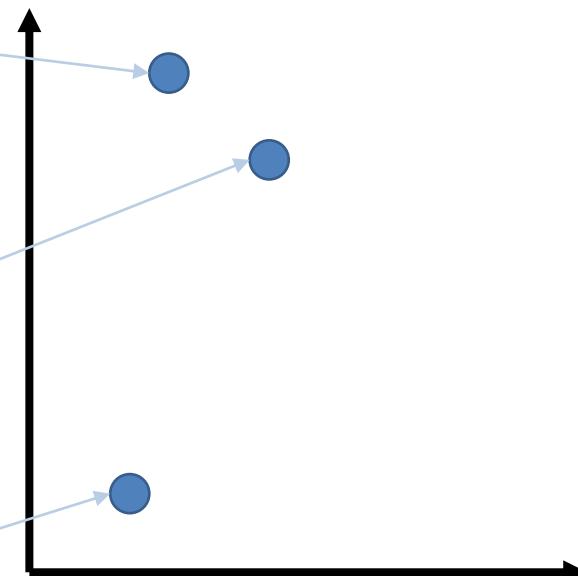


Projections:  
2-dimensional

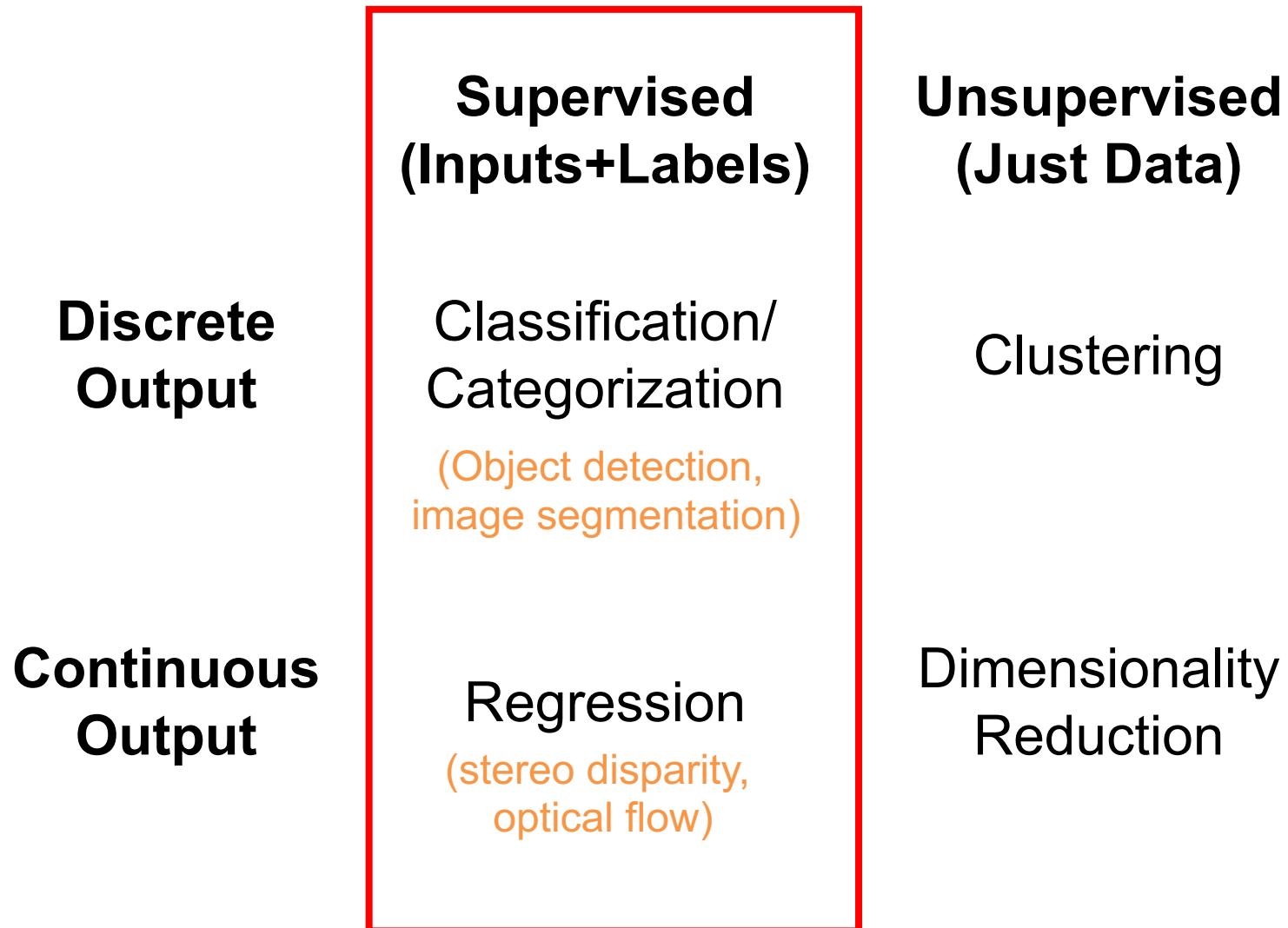
## Unsupervised Learning

**Data:**  $x$   
Just data, no labels!

**Goal:** Learn underlying  
*structure* in the data



# ML Problems in Vision



# **First Machine Learning Algorithm:**

# **Least Squares Linear Regression**

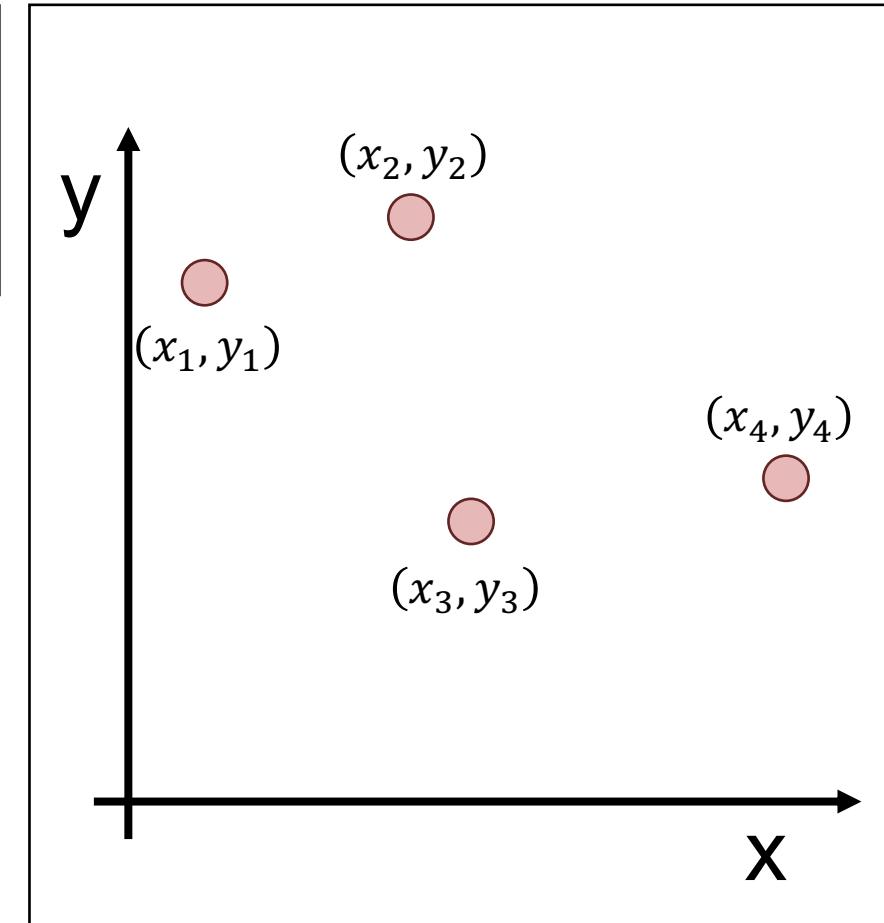
# Least Squares Linear Regression

# Least Squares Linear Regression

“Regression” = supervised learning with continuous outputs

Data:

$$(x_1, y_1), (x_2, y_2), \dots (x_N, y_N)$$
$$x_i, y_i \in \mathbb{R}$$



# Least Squares Linear Regression

“Linear” = Our model is a line

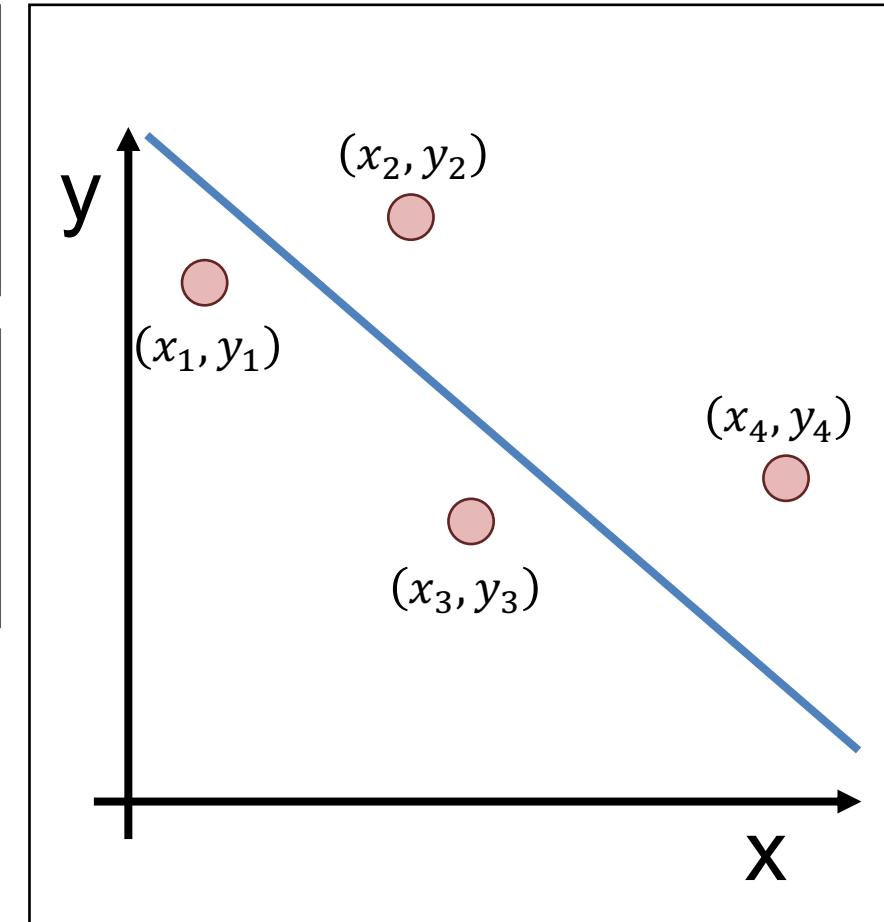
Data:

$$(x_1, y_1), (x_2, y_2), \dots (x_N, y_N)$$
$$x_i, y_i \in \mathbb{R}$$

Model:  $y = mx + b$

Or:  $\mathbf{x} = (x, 1)$ ;  $\mathbf{w} = (m, b)$

$$y = \mathbf{w} \cdot \mathbf{x}$$



# Least Squares Linear Regression

“Least squares” = Find the line that minimizes squared error

Data:

$$(x_1, y_1), (x_2, y_2), \dots (x_N, y_N)$$
$$x_i, y_i \in \mathbb{R}$$

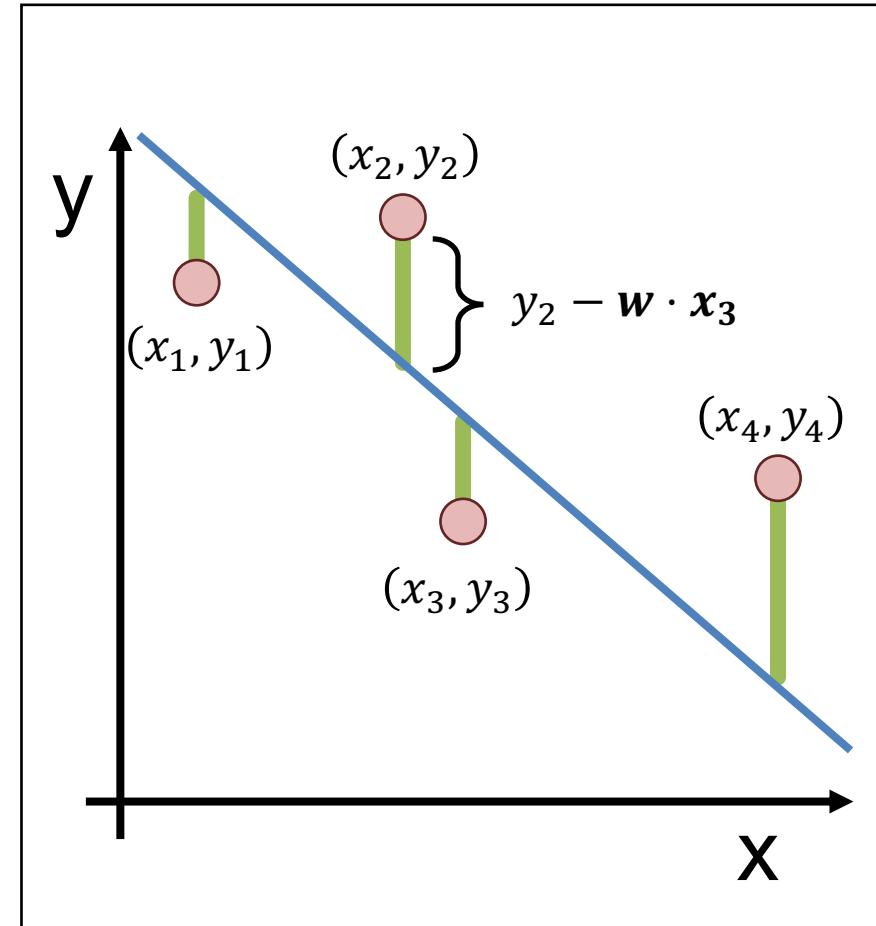
Model:  $y = mx + b$

Or:  $x = (x, 1)$ ;  $w = (m, b)$

$$y = w \cdot x$$

Training:

$$w^* = \arg \min_w \sum_{i=1}^N (y_i - w \cdot x_i)^2$$



# Solving Least Squares

$$\begin{aligned}\mathbf{w}^* &= \arg \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 \\ &= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2\end{aligned}$$

**Output:**

Vector of  
shape (N,)

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

**Inputs:**

Matrix of  
shape (N, 2)

$$\mathbf{X} = \begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_N & 1 \end{bmatrix}$$

**Weights:**

Vector of  
shape (2,)

$$\mathbf{w} = \begin{bmatrix} m \\ b \end{bmatrix}$$

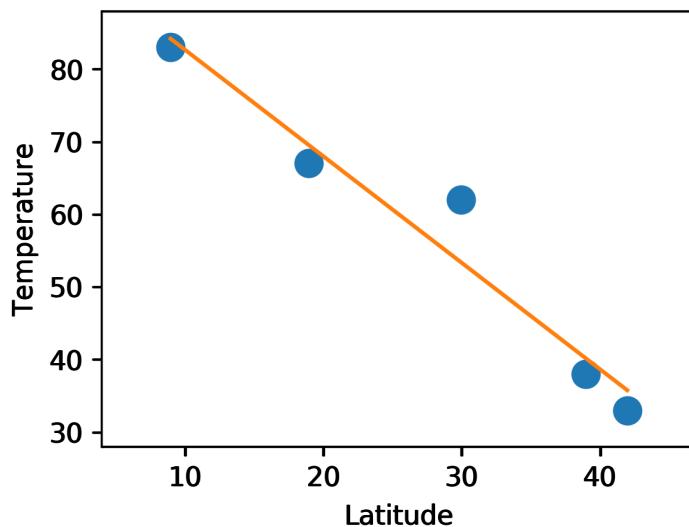
**Solution:**  $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

# Example: A Bad Weather Model

Given latitude, predict temperature by fitting a line

<u>City</u>	<u>Latitude (°)</u>	<u>Temp (F)</u>	<u>Training</u>
Ann Arbor	42	33	
Washington, DC	39	38	
Austin, TX	30	62	
Mexico City	19	67	
Panama City	9	83	

$$\mathbf{X}_{5 \times 2} = \begin{bmatrix} 42 & 1 \\ 39 & 1 \\ 30 & 1 \\ 19 & 1 \\ 9 & 1 \end{bmatrix} \quad \mathbf{y}_{5 \times 1} = \begin{bmatrix} 33 \\ 38 \\ 62 \\ 67 \\ 83 \end{bmatrix}$$



$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w}_{2 \times 1} = \begin{bmatrix} -1.47 \\ 97 \end{bmatrix}$$

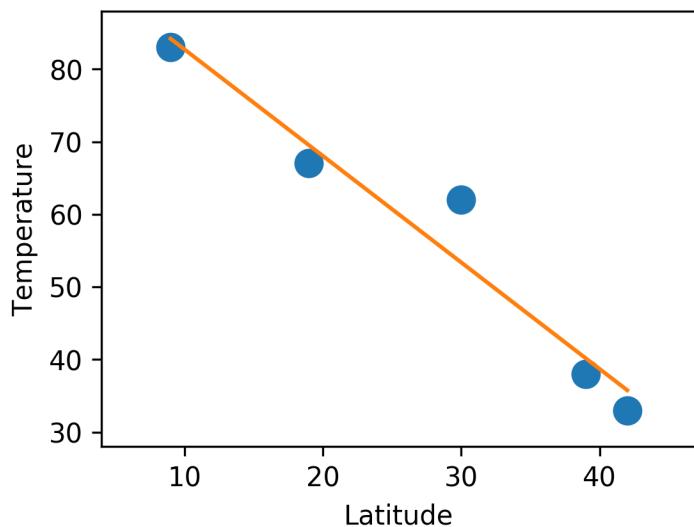
$$\text{Temp} = -1.47 * \text{Lat} + 97$$

# Example: A Bad Weather Model

Given latitude, predict temperature by fitting a line

<u>City</u>	<u>Latitude (°)</u>	<u>Temp (F)</u>	<u>Prediction</u>	<u>Error</u>	
Ann Arbor	42	33	35.3	2.3	
Washington, DC	39	38	39.7	1.7	
Austin, TX	30	62	52.9	10.9	
Mexico City	19	67	69.1	2.1	
Panama City	9	83	83.8	0.8	

Seems  
good!

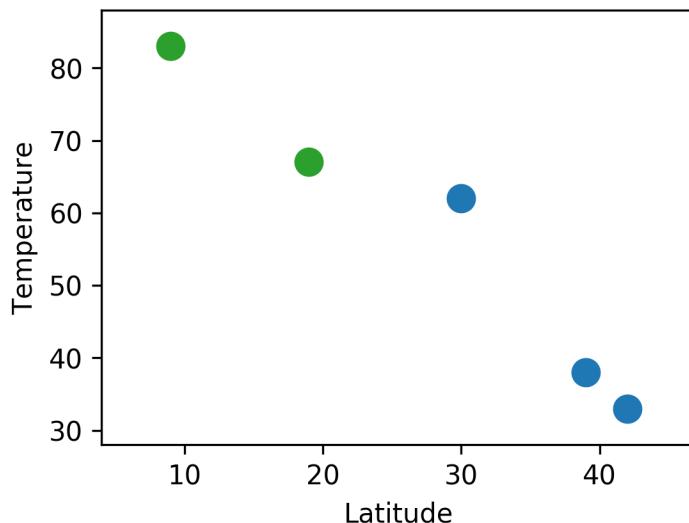


**Problem:** In ML we don't care about training set performance; we want models that **generalize** to new data

# Example: A Bad Weather Model

Given latitude, predict temperature by fitting a line

	<u>City</u>	<u>Latitude (°)</u>	<u>Temp (F)</u>
<b>Train</b>	Ann Arbor	42	33
	Washington, DC	39	38
	Austin, TX	30	62
<b>Test</b>	Mexico City	19	67
	Panama City	9	83



**Problem:** In ML we don't care about training set performance; we want models that **generalize** to new data

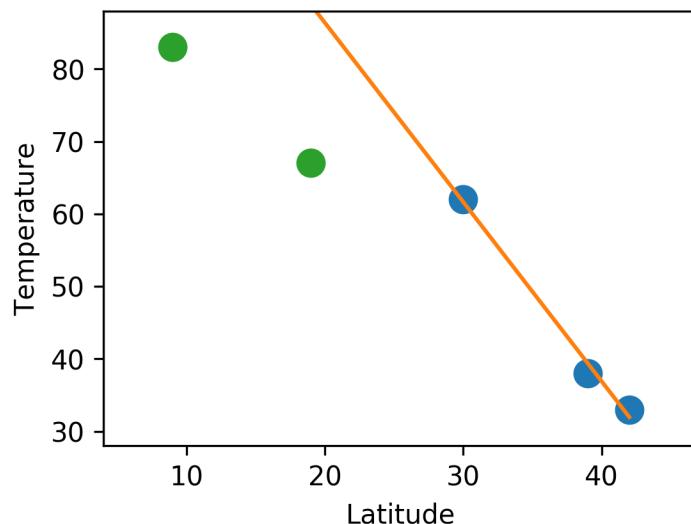
**Solution:** Split dataset into **train** and **test** sets

# Example: A Bad Weather Model

Given latitude, predict temperature by fitting a line

	<u>City</u>	<u>Latitude (°)</u>	<u>Temp (F)</u>	<u>Prediction</u>	<u>Error</u>
Train	Ann Arbor	42	33	31.9	1.0
	Washington, DC	39	38	39.4	1.4
	Austin, TX	30	62	61.7	0.3
Test	Mexico City	19	67	88.9	21.9
	Panama City	9	83	113.6	30.6

**Problem:** Low error on train, high error on test

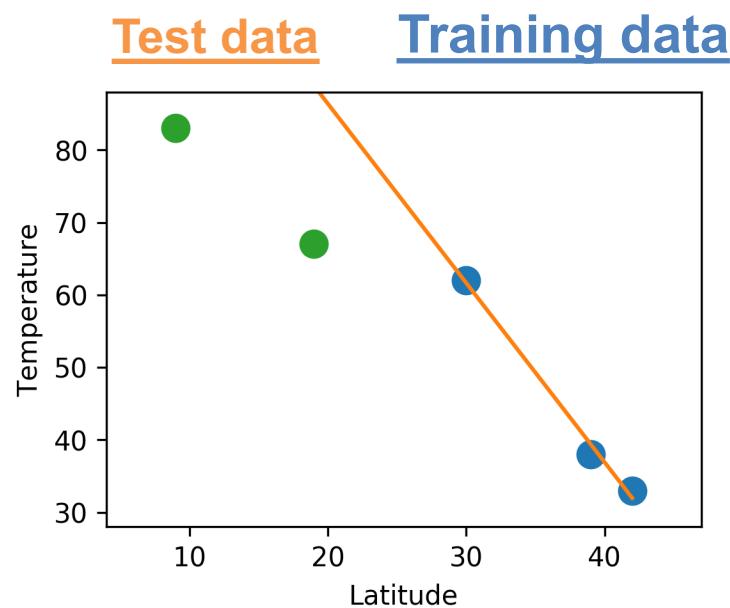


**Problem:** In ML we don't care about training set performance; we want models that **generalize** to new data

**Solution:** Split dataset into train and test sets  
Fit on train set, evaluate on both

# Overfitting,

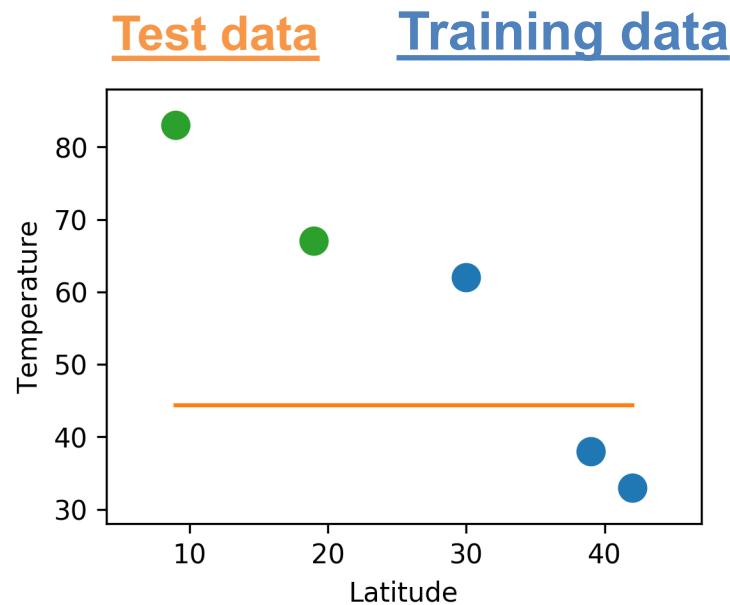
**Overfitting:** Model fits noise in the training set  
and doesn't generalize well to new data  
Model is too flexible



# Overfitting, Underfitting,

**Overfitting**: Model fits noise in the training set  
and doesn't generalize well to new data  
Model is too flexible

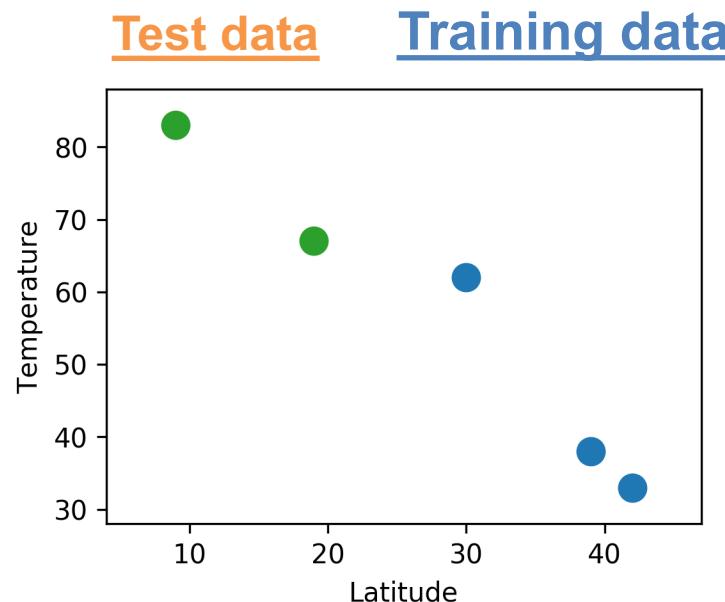
**Underfitting**: Model doesn't fit the training data  
well, high error on both train and test  
Model is not flexible enough



# Overfitting, Underfitting, Regularization

**Overfitting:** Model fits noise in the training set and doesn't generalize well to new data  
Model is too flexible

**Underfitting:** Model doesn't fit the training data well, high error on both train and test  
Model is not flexible enough



**Regularization:** Penalize model complexity to prevent overfitting and improve generalization

## Least Squares

$$\arg \min_w \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

## L2-Regularized Least Squares

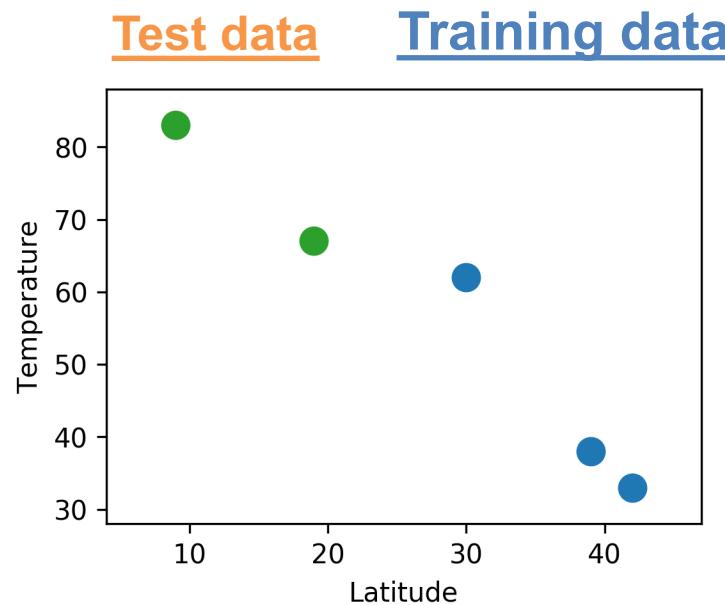
$$\arg \min_w \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

# Overfitting, Underfitting, Regularization

**Overfitting:** Model fits noise in the training set and doesn't generalize well to new data  
Model is too flexible

**Underfitting:** Model doesn't fit the training data well, high error on both train and test  
Model is not flexible enough



**Regularization:** Penalize model complexity to prevent overfitting and improve generalization

Least Squares

$$\arg \min_w \|\mathbf{y} - \mathbf{Xw}\|^2$$

L2-Regularized Least Squares

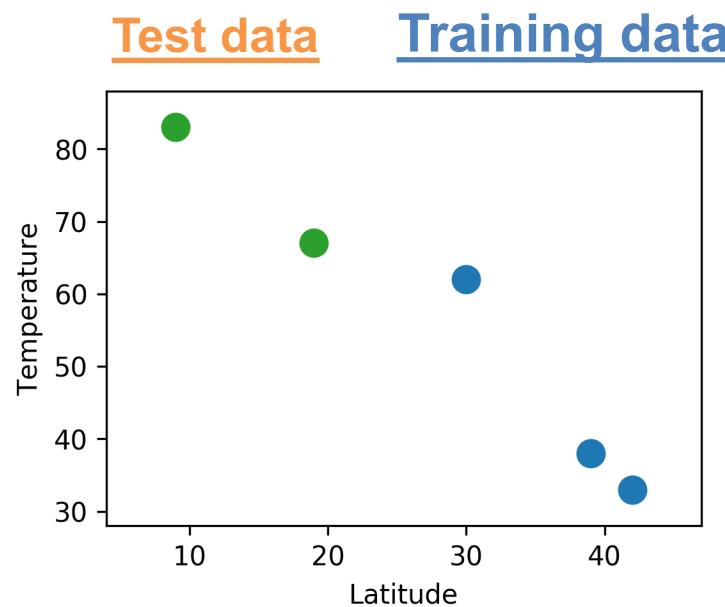
$$\arg \min_w \|\mathbf{y} - \mathbf{Xw}\|^2 + \lambda \|\mathbf{w}\|^2$$

Fit training data

# Overfitting, Underfitting, Regularization

**Overfitting:** Model fits noise in the training set and doesn't generalize well to new data  
Model is too flexible

**Underfitting:** Model doesn't fit the training data well, high error on both train and test  
Model is not flexible enough



**Regularization:** Penalize model complexity to prevent overfitting and improve generalization

Least Squares

$$\arg \min_w \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

L2-Regularized Least Squares

$$\arg \min_w \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

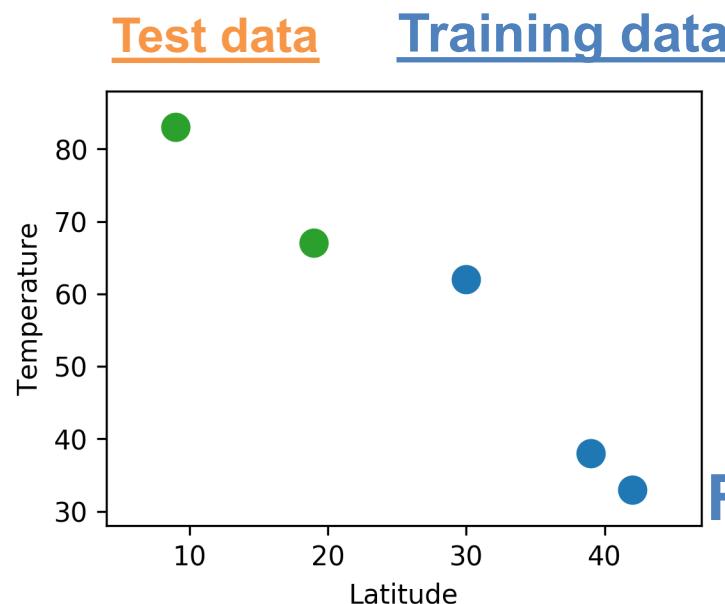
Fit training data

↑  
Penalize complexity

# Overfitting, Underfitting, Regularization

**Overfitting:** Model fits noise in the training set and doesn't generalize well to new data  
Model is too flexible

**Underfitting:** Model doesn't fit the training data well, high error on both train and test  
Model is not flexible enough



**Regularization:** Penalize model complexity to prevent overfitting and improve generalization

## Least Squares

$$\arg \min_w \|\mathbf{y} - \mathbf{Xw}\|^2$$

## L2-Regularized Least Squares

$$\arg \min_w \|\mathbf{y} - \mathbf{Xw}\|^2 + \lambda \|\mathbf{w}\|^2$$

Fit training data

Regularization Strength

Penalize complexity

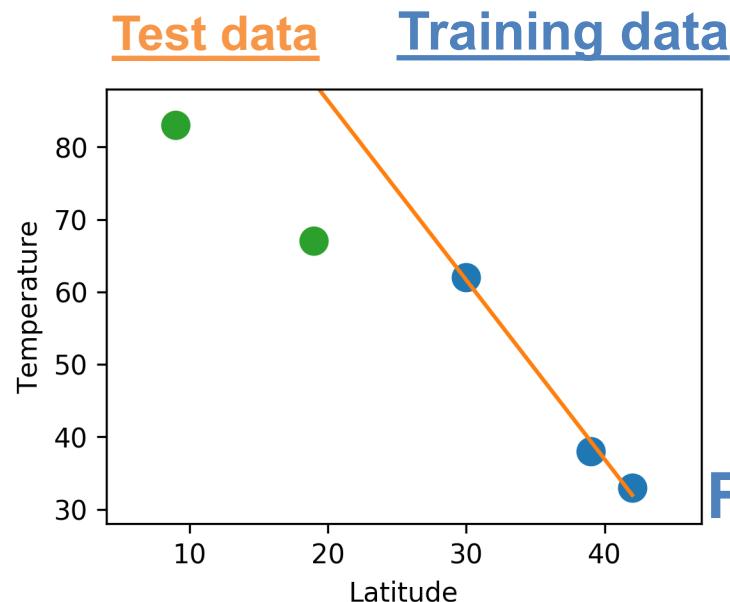
# Overfitting, Underfitting, Regularization

**Overfitting:** Model fits noise in the training set and doesn't generalize well to new data  
Model is too flexible

**Underfitting:** Model doesn't fit the training data well, high error on both train and test  
Model is not flexible enough

**Regularization:** Penalize model complexity to prevent overfitting and improve generalization

$\lambda = 0.00$ : No regularization; model overfits



L2-Regularized Least Squares  
$$\arg \min_w \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

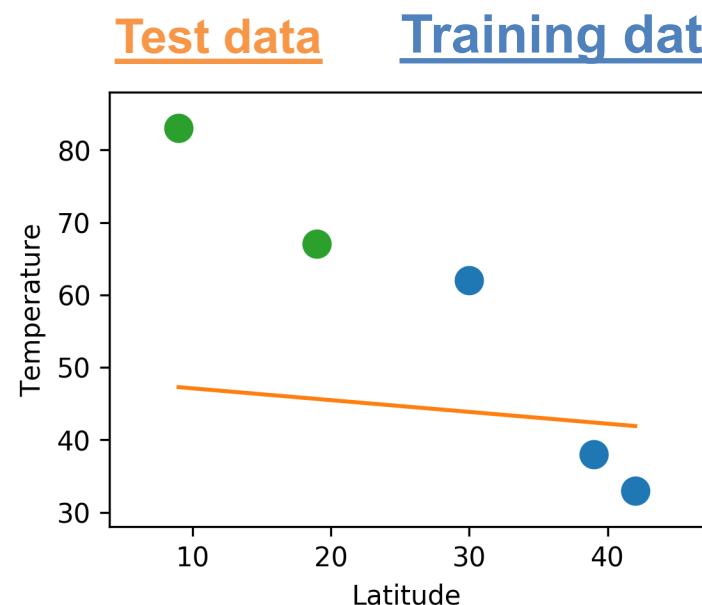
Fit training data      Regularization Strength      Penalize complexity

# Overfitting, Underfitting, Regularization

**Overfitting:** Model fits noise in the training set and doesn't generalize well to new data  
Model is too flexible

**Underfitting:** Model doesn't fit the training data well, high error on both train and test  
Model is not flexible enough

**Regularization:** Penalize model complexity to prevent overfitting and improve generalization



$\lambda = 0.00$ : No regularization; model overfits

$\lambda = 0.10$ : Too much regularization, underfitting

L2-Regularized Least Squares

$$\arg \min_w \|y - Xw\|^2 + \lambda \|w\|^2$$

Fit training data

Regularization Strength

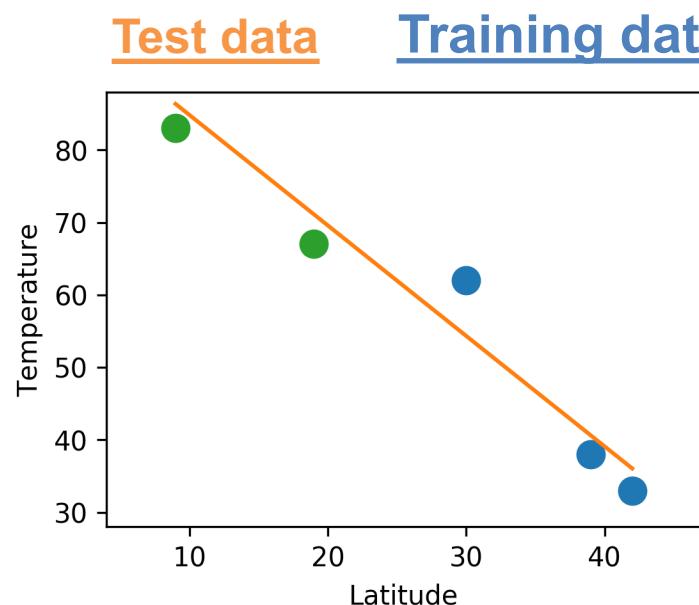
Penalize complexity

# Overfitting, Underfitting, Regularization

**Overfitting:** Model fits noise in the training set and doesn't generalize well to new data  
Model is too flexible

**Underfitting:** Model doesn't fit the training data well, high error on both train and test  
Model is not flexible enough

**Regularization:** Penalize model complexity to prevent overfitting and improve generalization



$\lambda = 0.00$ : No regularization; model overfits

$\lambda = 0.10$ : Too much regularization, underfitting

$\lambda = 0.02$ : Just right! Good fit and generalization

L2-Regularized Least Squares

$$\arg \min_w \|y - Xw\|^2 + \lambda \|w\|^2$$

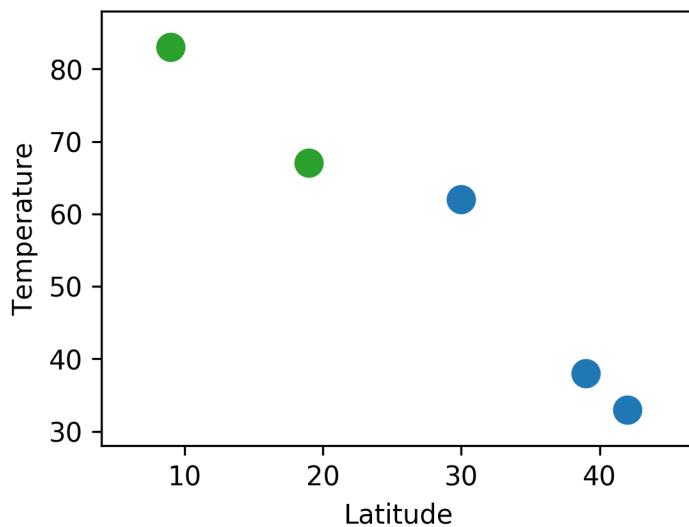
Fit training data

Regularization Strength

Penalize complexity

# Regularization

	<u>City</u>	<u>Latitude (°)</u>	<u>Temp (F)</u>
<b>Train</b>	Ann Arbor	42	33
	Washington, DC	39	38
	Austin, TX	30	62
<b>Test</b>	Mexico City	19	67
	Panama City	9	83

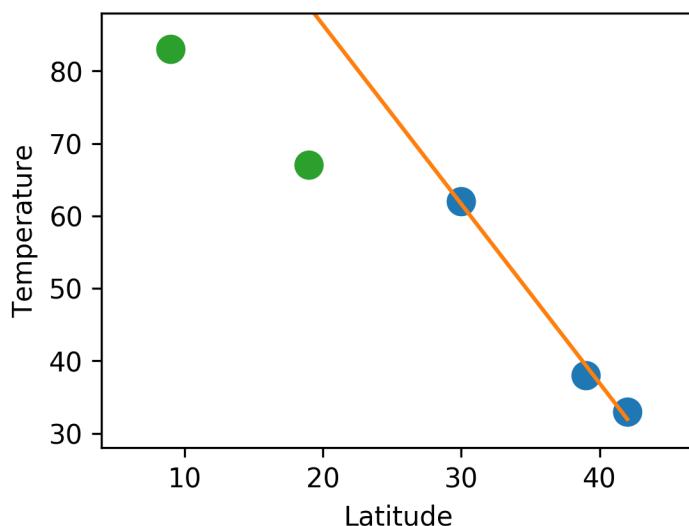


L2-Regularized Least Squares  
$$\arg \min_w \|y - Xw\|^2 + \lambda \|w\|^2$$

Fit training data      Regularization Strength      Penalize complexity

# Regularization

		No regularization ( $\lambda = 0$ )		
	<u>City</u>	<u>Latitude (°)</u>	<u>Temp (F)</u>	<u>Prediction</u>
<b>Train</b>	Ann Arbor	42	33	31.9
	Washington, DC	39	38	39.4
	Austin, TX	30	62	61.7
<b>Test</b>	Mexico City	19	67	88.9
	Panama City	9	83	113.6



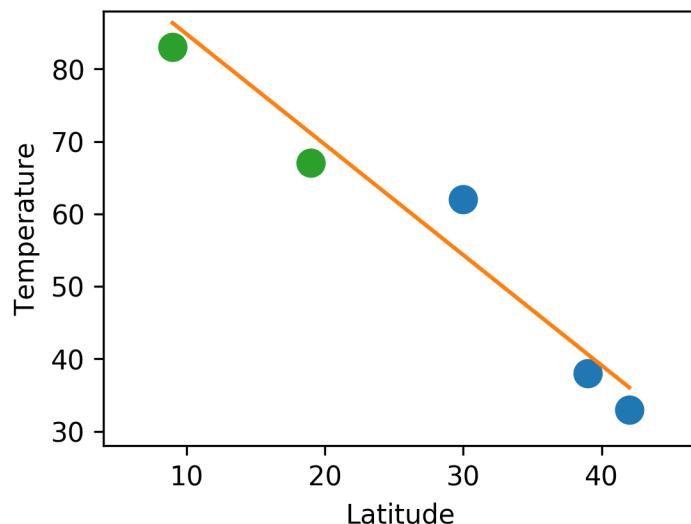
L2-Regularized Least Squares

$$\arg \min_w \|y - Xw\|^2 + \lambda \|w\|^2$$

Fit training data      Regularization Strength      Penalize complexity

# Regularization

	<u>City</u>	<u>Latitude (°)</u>	<u>Temp (F)</u>	No regularization ( $\lambda = 0$ )		Regularized ( $\lambda = 0.02$ )	
				<u>Prediction</u>	<u>Error</u>	<u>Prediction</u>	<u>Error</u>
<b>Train</b>	Ann Arbor	42	33	31.9	1.0	36.0	3.0
	Washington, DC	39	38	39.4	1.4	40.6	2.6
	Austin, TX	30	62	61.7	0.3	54.3	7.7
<b>Test</b>	Mexico City	19	67	88.9	21.9	71.1	4.1
	Panama City	9	83	113.6	30.6	86.4	3.4



L2-Regularized Least Squares

$$\arg \min_w \|y - Xw\|^2 + \lambda \|w\|^2$$

Fit training data      Regularization Strength      Penalize complexity

# Parameters and Hyperparameters

L2-Regularized Least Squares

$$\arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

**Parameter ( $w$ ):** **(Automatically)** selected during training by fitting to training data

**Hyperparameter ( $\lambda$ ):** **(Manually)** chosen before training, does not depend on training data

**Question:** How to choose hyperparameters?

# Choosing Hyperparameters

**Idea #1:** Choose hyperparameters that work best on the data

**BAD:**  $\lambda = 0$  always works best on training data

Your Dataset

**Idea #2:** Split data into **train** and **test**, choose hyperparameters that work best on test data

**BAD:** No idea how we will perform on new data

train

test

**Idea #3:** Split data into **train**, **val**, and **test**; choose hyperparameters on val and evaluate on test

**Better**  
!

train

validation

test

# Choosing Hyperparameters

Your Dataset

**Idea #4: Cross-Validation:** Split data into **folds**,  
try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 4	fold 5	test
--------	--------	--------	--------	--------	------

fold 1	fold 2	fold 3	fold 4	fold 5	test
--------	--------	--------	--------	--------	------

fold 1	fold 2	fold 3	fold 4	fold 5	test
--------	--------	--------	--------	--------	------

Useful for small datasets, but (unfortunately) not  
used too frequently in deep learning

# Image Classification: Core Vision Task

**Input:** image



This image by [Nikita](#) is  
licensed under [CC-BY 2.0](#)

**Output:** Assign image to  
one of a fixed set of  
categories

cat

bird

deer

dog

truck



# Classification with Least Squares

$x_i \in \mathbb{R}^D$  is image feature

$y_i \in \mathbb{R}^C$  is **one-hot** label

$y_{i,c} = 1$  if  $x_i$  has category c, 0 otherwise

Training ( $\mathbf{x}_i, \mathbf{y}_i$ ):

$$\arg \min_W \sum_{i=1}^n \|\mathbf{W}\mathbf{x}_i - \mathbf{y}_i\|^2$$

Inference ( $\mathbf{x}$ ):

$$\mathbf{W}\mathbf{x} > t$$

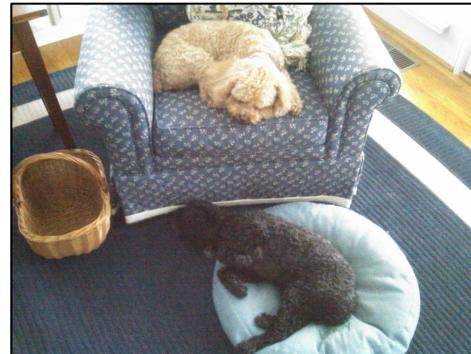
Unprincipled in theory, but often effective in practice

The reverse (regression via discrete bins) is also common

Rifkin, Yeo, Poggio. *Regularized Least Squares Classification* (<http://cbcl.mit.edu/publications/ps/rlsc.pdf>). 2003  
Redmon, Divvala, Girshick, Farhadi. *You Only Look Once: Unified, Real-Time Object Detection*. CVPR 2016.

# Linear Classifiers

$$y_i = \mathbf{w}x_i + b, \quad x_i \in \mathbb{R}^D \text{ is image feature}$$



Model – one weight per class:  $\mathbf{w}_0, \mathbf{w}_1, \mathbf{w}_2$

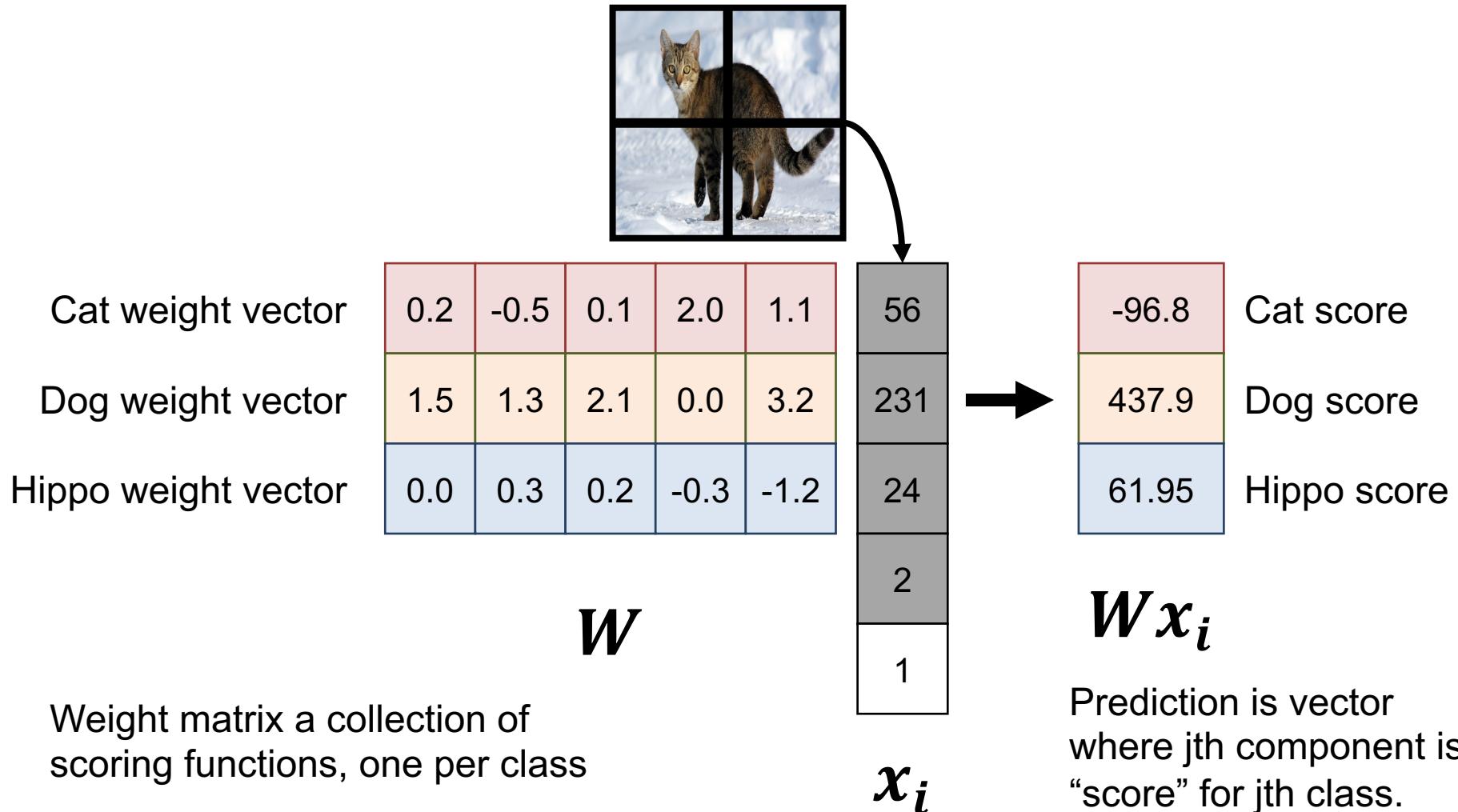
$\mathbf{w}_0^T \mathbf{x}$  big if cat

$\mathbf{w}_1^T \mathbf{x}$  big if dog

$\mathbf{w}_2^T \mathbf{x}$  big if hippo

Stack together:  $W_{3xD}$

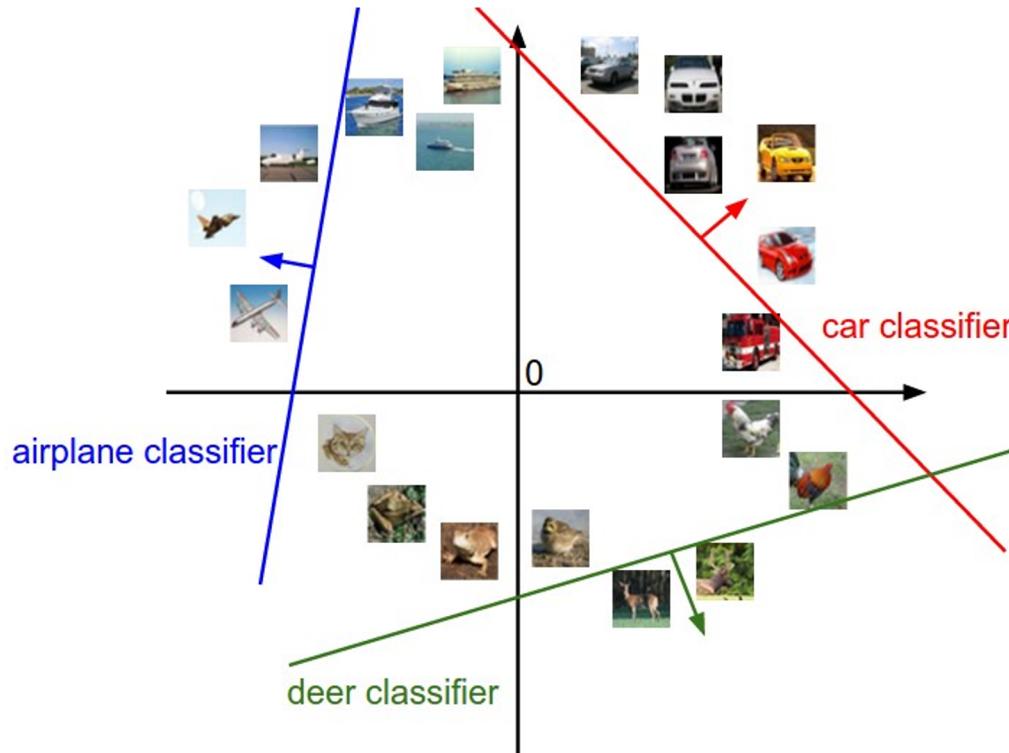
# Linear Classifiers



Weight matrix a collection of scoring functions, one per class

# Linear Classifiers: Geometric Intuition

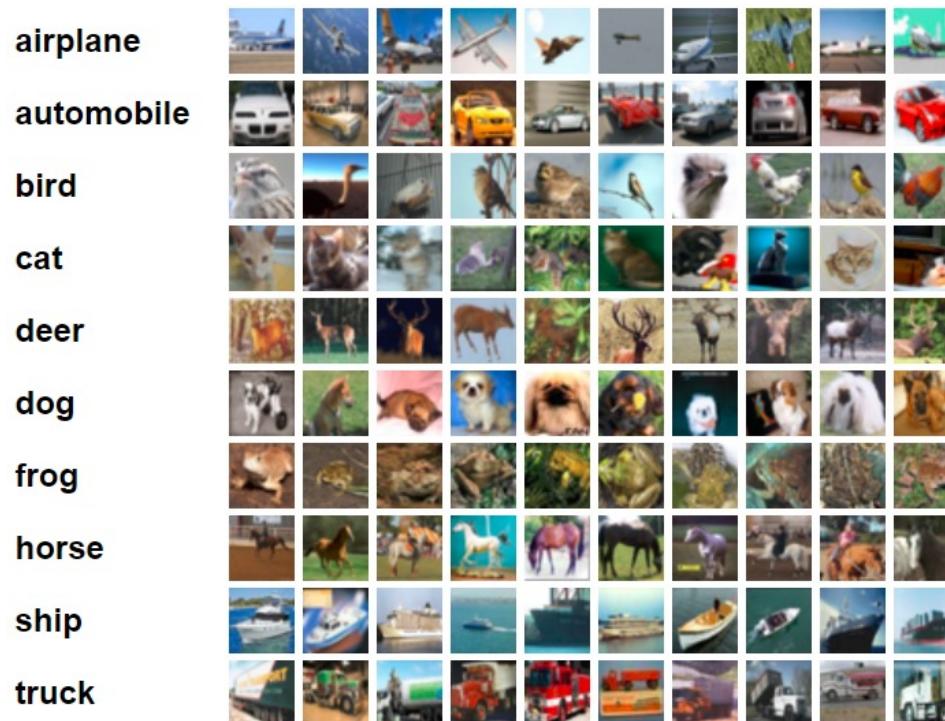
What does a linear classifier look like in 2D?



**Be aware:** Intuition from 2D doesn't always carry over into high-dimensional spaces. See: *On the Surprising Behavior of Distance Metrics in High Dimensional Space*. Charu, Hinneburg, Keim. ICDT 2001

# Linear Classifiers: Visual Intuition

CIFAR 10:  
32x32x3 Images, 10 Classes



- Turn each image into feature by unrolling all pixels
- Train a linear model to recognize 10 classes

# Linear Classifiers: Visual Intuition

Decision rule is  $\mathbf{w}^T \mathbf{x}$ . If  $w_i$  is big, then  
big values of  $x_i$  are indicative of the  
class.

## Deer or Plane?



# Linear Classifiers: Visual Intuition

Decision rule is  $\mathbf{w}^T \mathbf{x}$ . If  $w_i$  is big, then  
big values of  $x_i$  are indicative of the  
class.



# So Far: Linear Score Function



Model – one weight per class:  $w_0, w_1, w_2$

$$\begin{aligned} w_0^T x & \text{ big if cat} \\ w_1^T x & \text{ big if dog} \\ w_2^T x & \text{ big if hippo} \end{aligned}$$

The diagram consists of three horizontal arrows pointing upwards from the equations to the corresponding weights. The first arrow is red, the second is orange, and the third is blue. They all point to the right, aligning with the vertical position of each weight.

Stack together:  $W_{3x_F}$  where  $x$  is in  $R^F$

How do we know which  $W$  is best?

# Choosing W: Loss Function

A **loss function** tells how good our current classifier is

Low loss = good classifier  
High loss = bad classifier

(Also called: **objective function**; **cost function**)

Given a dataset

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

of images  $x_i$  and labels  $y_i$ ,

Loss for a single example is:  
 $L_i(f(x_i, W), y_i)$

Loss for the dataset is

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

# Cross-Entropy Loss

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

Classifier scores

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



cat      **3.2**

car      5.1

frog     -1.7

# Cross-Entropy Loss

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

Classifier scores  
 $s = f(x_i, W)$



Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

cat      **3.2**

car      5.1

frog     -1.7

# Cross-Entropy Loss

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

Classifier scores  
 $s = f(x_i, W)$



Softmax function

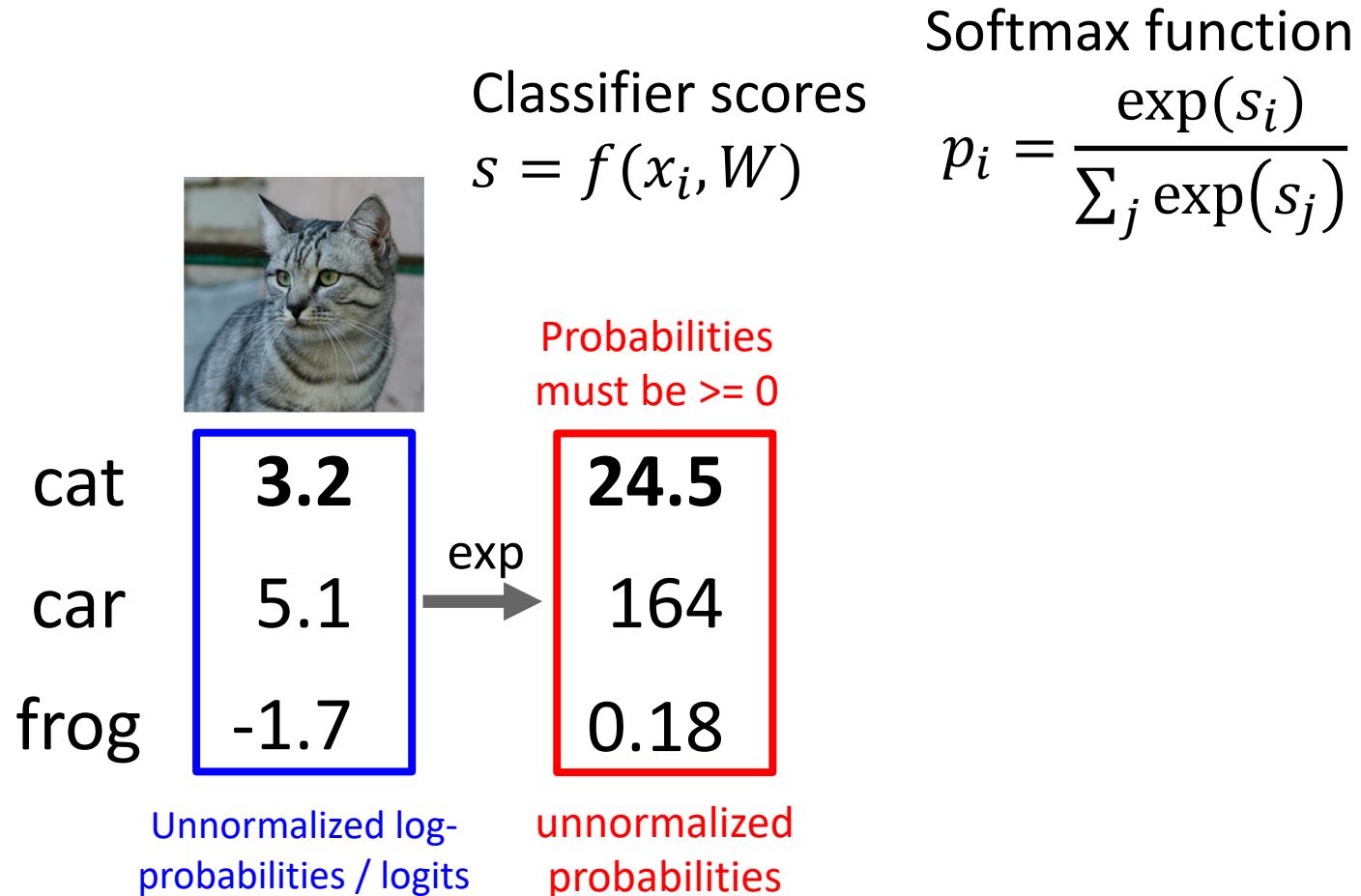
$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

cat	3.2
car	5.1
frog	-1.7

Unnormalized log-  
probabilities / logits

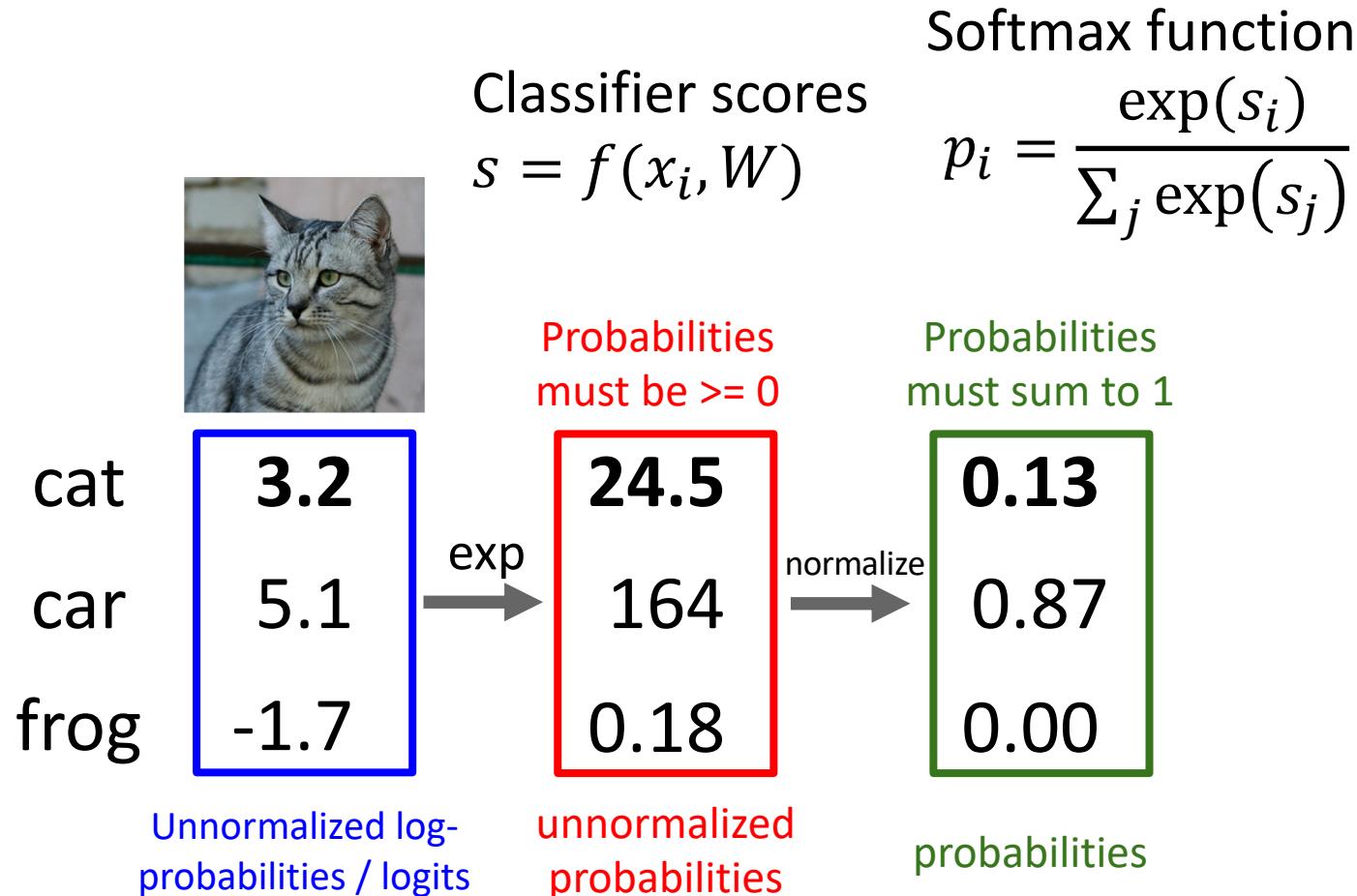
# Cross-Entropy Loss

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



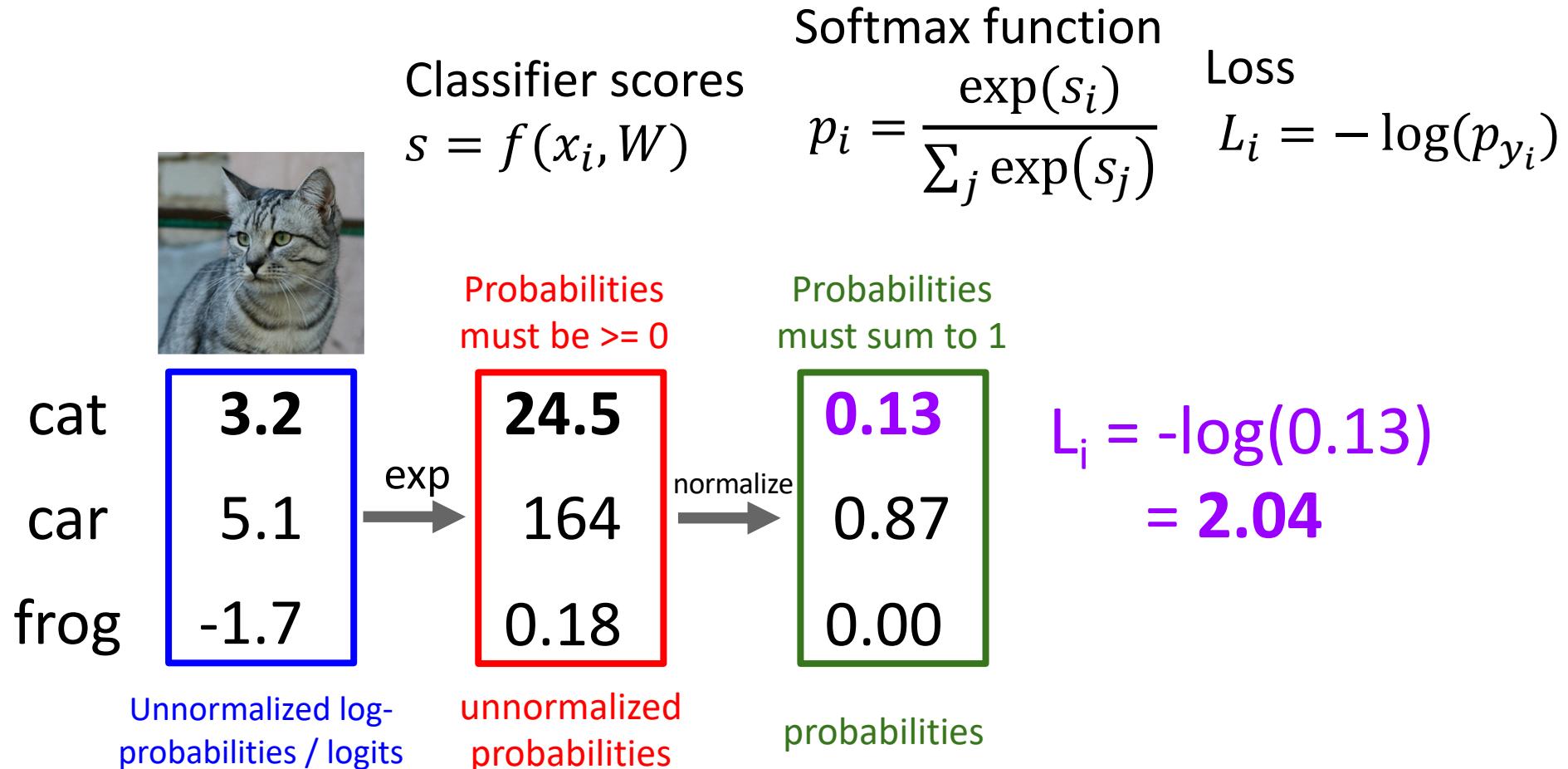
# Cross-Entropy Loss

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



# Cross-Entropy Loss

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



# Cross-Entropy Loss

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

	
cat	3.2
car	5.1
frog	-1.7

Unnormalized log-probabilities / logits

Classifier scores  
 $s = f(x_i, W)$

Probabilities must be  $\geq 0$

24.5
164
0.18

$\exp$

unnormalized probabilities

Softmax function

$$p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

Loss

$$L_i = -\log(p_{y_i})$$

Probabilities must sum to 1

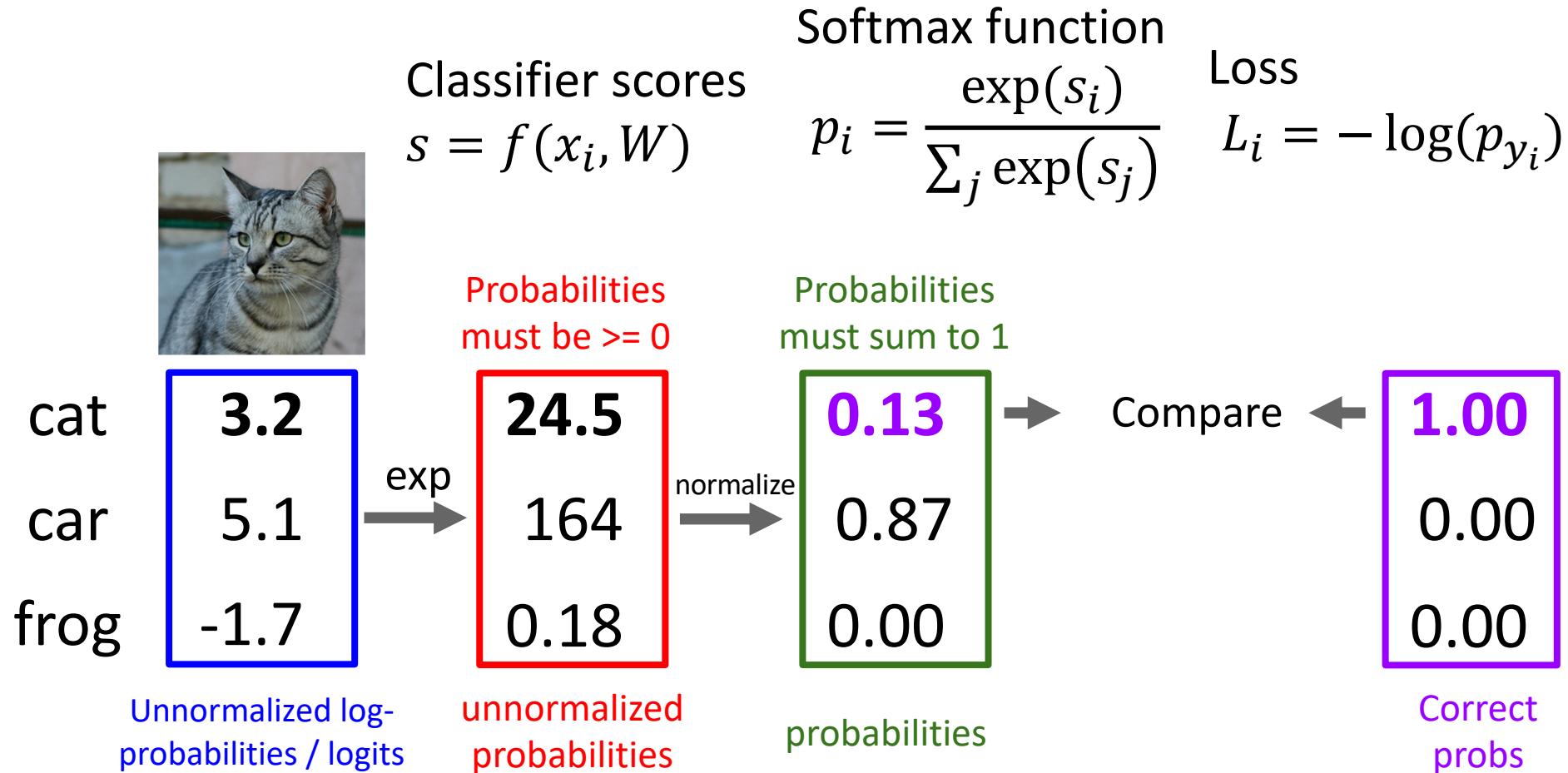
0.13
0.87
0.00

$$L_i = -\log(0.13) \\ = 2.04$$

**Maximum Likelihood Estimation**  
Choose weights to maximize the likelihood of the observed data

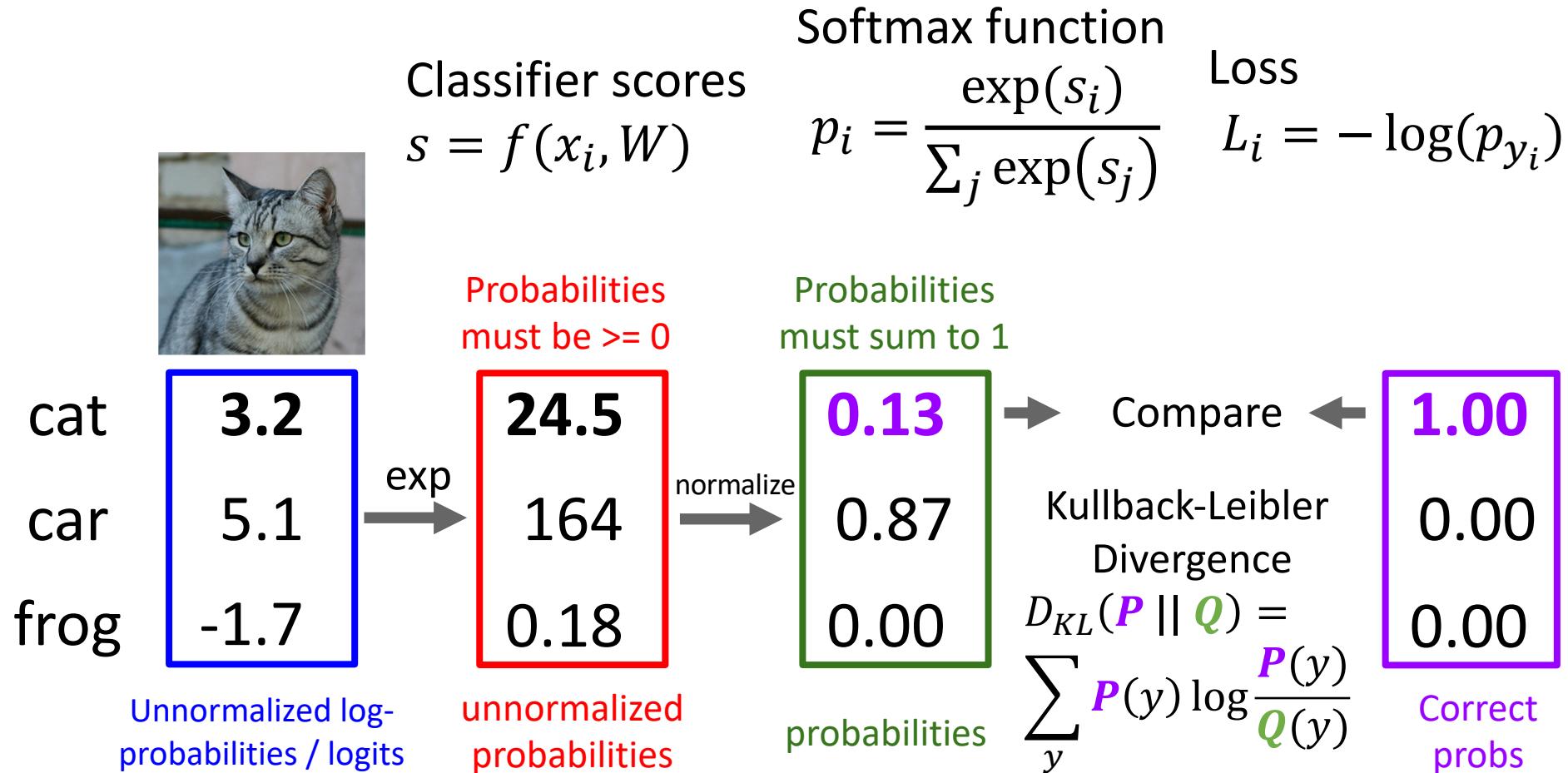
# Cross-Entropy Loss

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



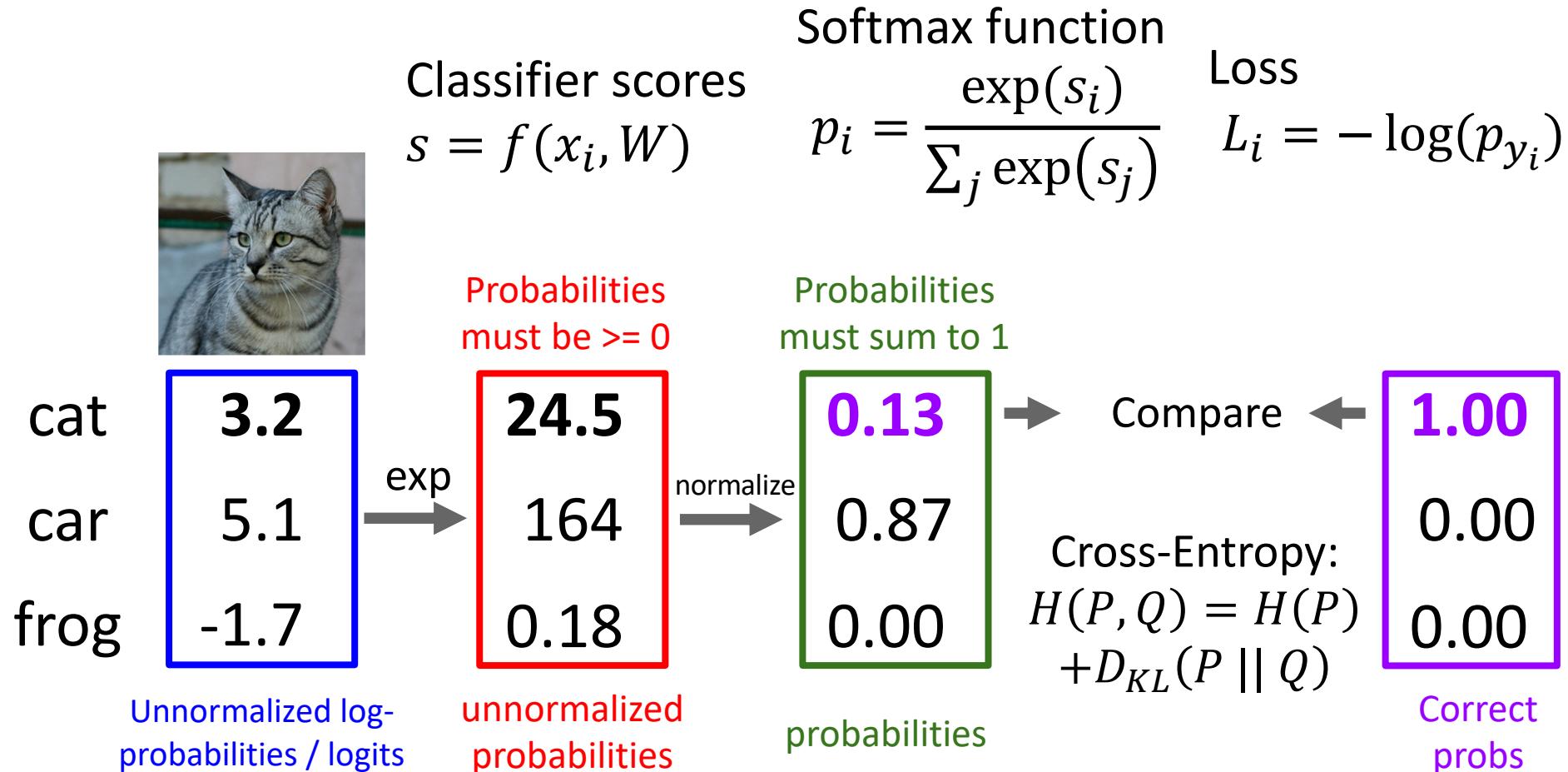
# Cross-Entropy Loss

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



# Cross-Entropy Loss

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



# Cross-Entropy Loss

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



cat	<b>3.2</b>
car	5.1
frog	-1.7

Classifier scores  
 $s = f(x_i, W)$

Softmax function  
 $p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$

Loss  
 $L_i = -\log(p_{y_i})$

Putting it all together:

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

# Cross-Entropy Loss

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



cat	<b>3.2</b>
car	5.1
frog	-1.7

Classifier scores  
 $s = f(x_i, W)$

Softmax function  
 $p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$

Loss  
 $L_i = -\log(p_{y_i})$

Putting it all together:

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

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

# Cross-Entropy Loss

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



cat	<b>3.2</b>
car	5.1
frog	-1.7

Classifier scores  
 $s = f(x_i, W)$

Softmax function  
 $p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}$

Loss  
 $L_i = -\log(p_{y_i})$

Putting it all together:

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

**Q:** If all scores are small random values, what is the loss?

# Next Class

- A bit more about loss function for classification
  - SVM (hinge) loss
- Optimization
  - How to choose the optimal  $W$  ( $W^*$ )