# Machine learning

370: Intro to Computer Vision

Subhransu Maji

April 15 & 17, 2025

# Image classification



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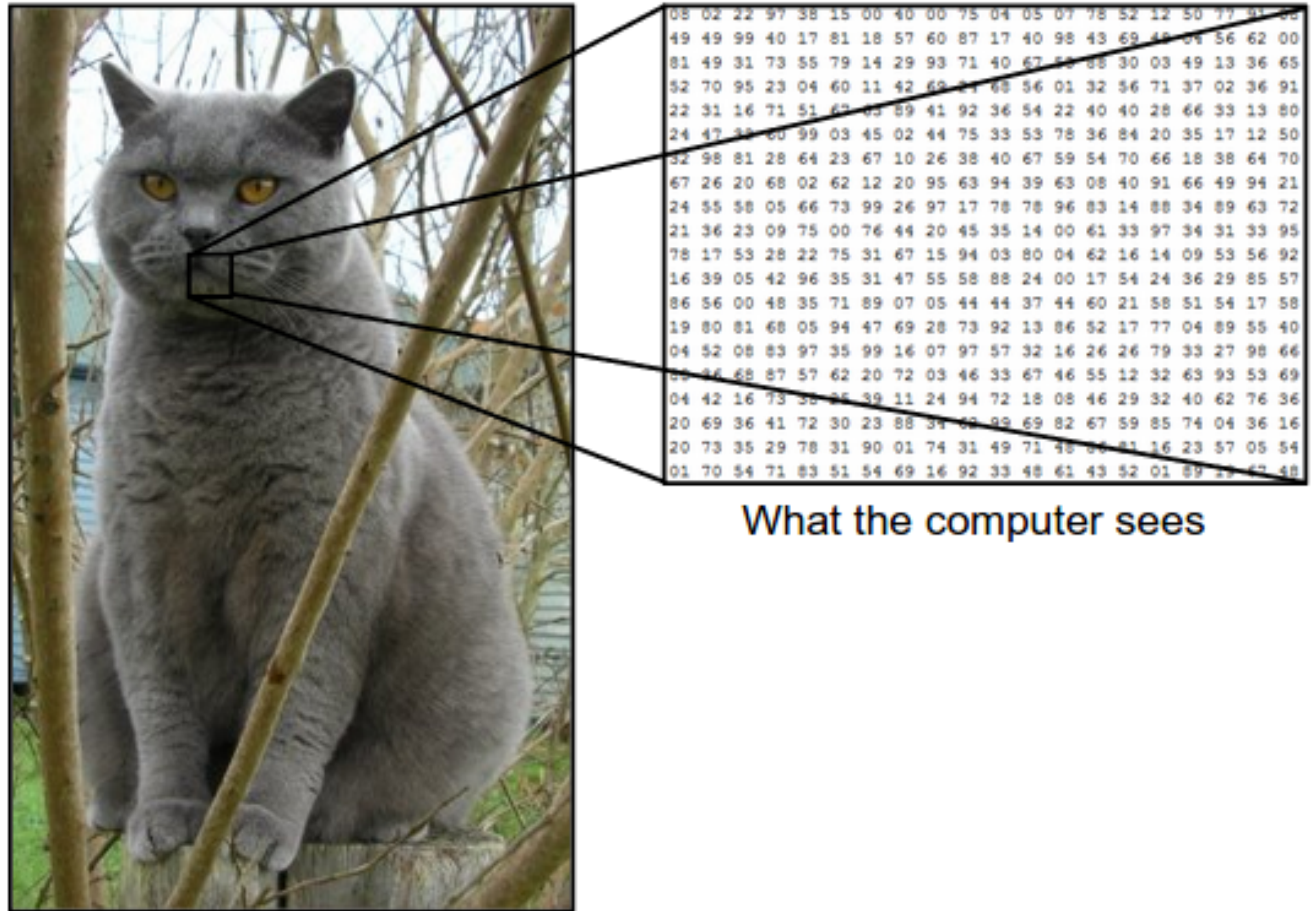
$\longrightarrow$     cat

# Challenges: Semantic gap

Images are represented as 3D arrays of numbers, with integers between [0, 255].
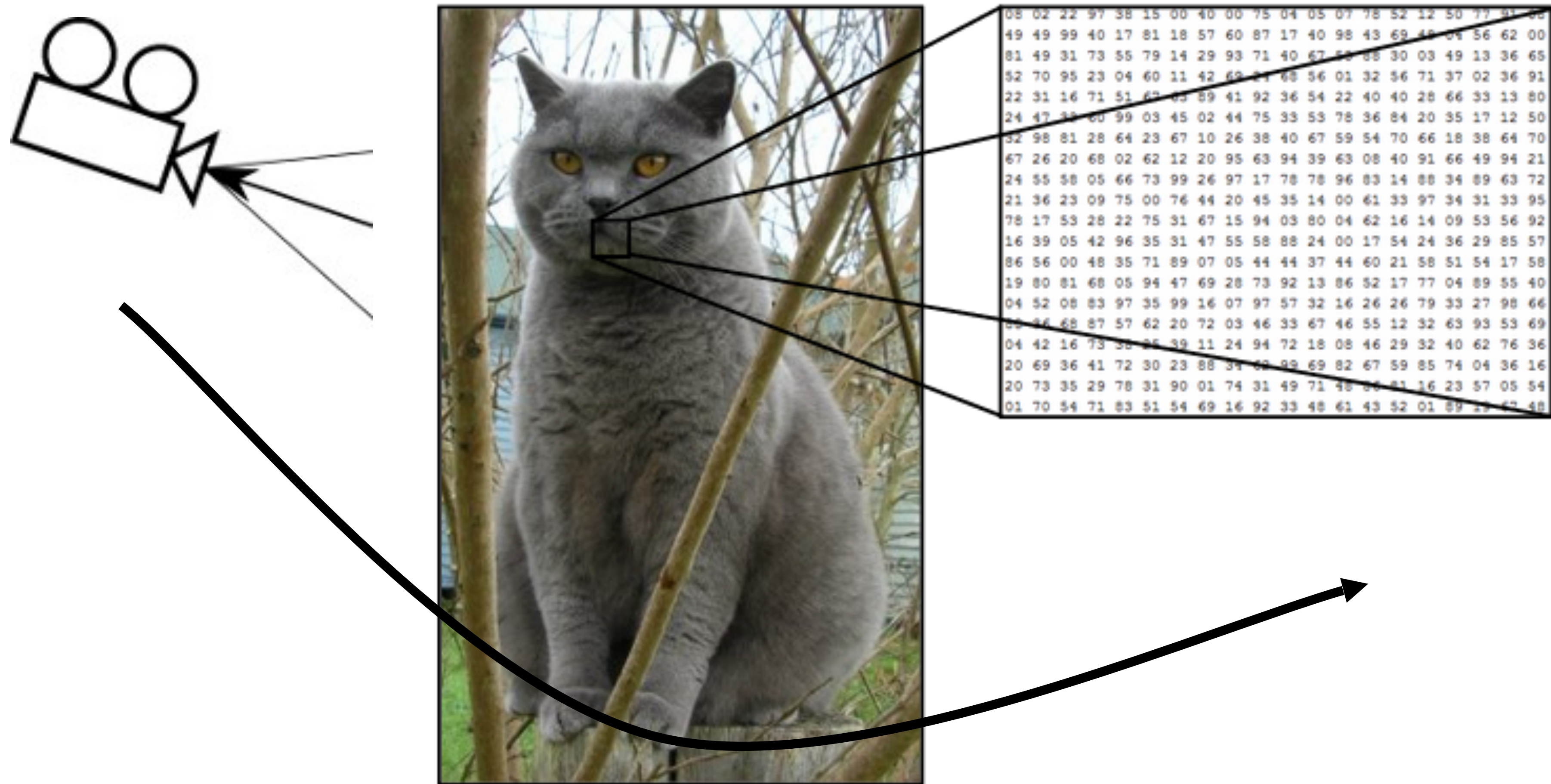
E.g.
300 x 100 x 3

(3 for 3 color channels RGB)



What the computer sees
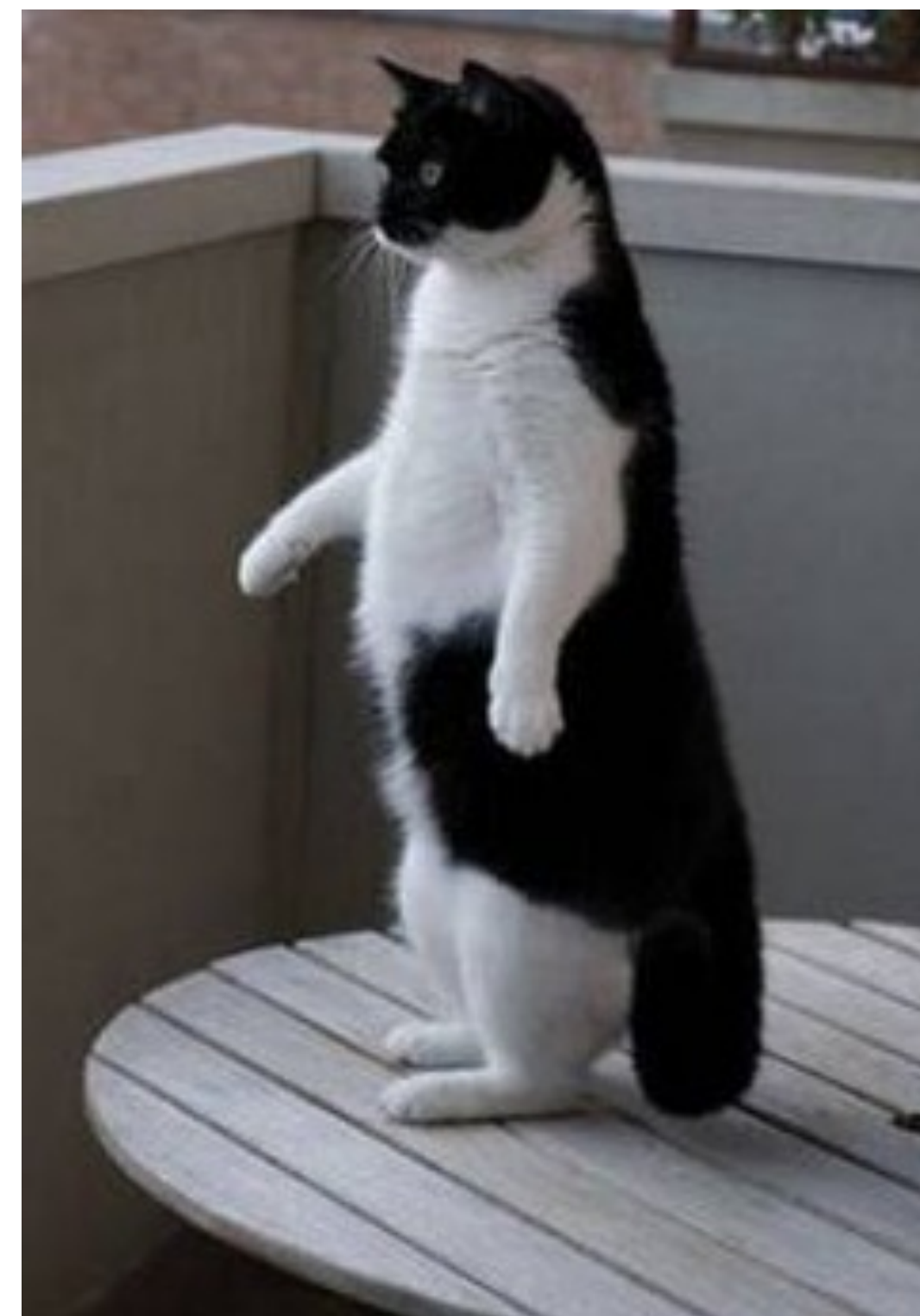
# Challenges: Viewpoint Variation

# Challenges: Illumination

# Challenges: Deformation

# Challenges: Occlusion

# Challenges: Background clutter

# Challenges: Intraclass variation

# Writing an image classifier

```python
def predict(image):
    # ????
    return class_label
```

Unlike e.g. sorting a list of numbers,

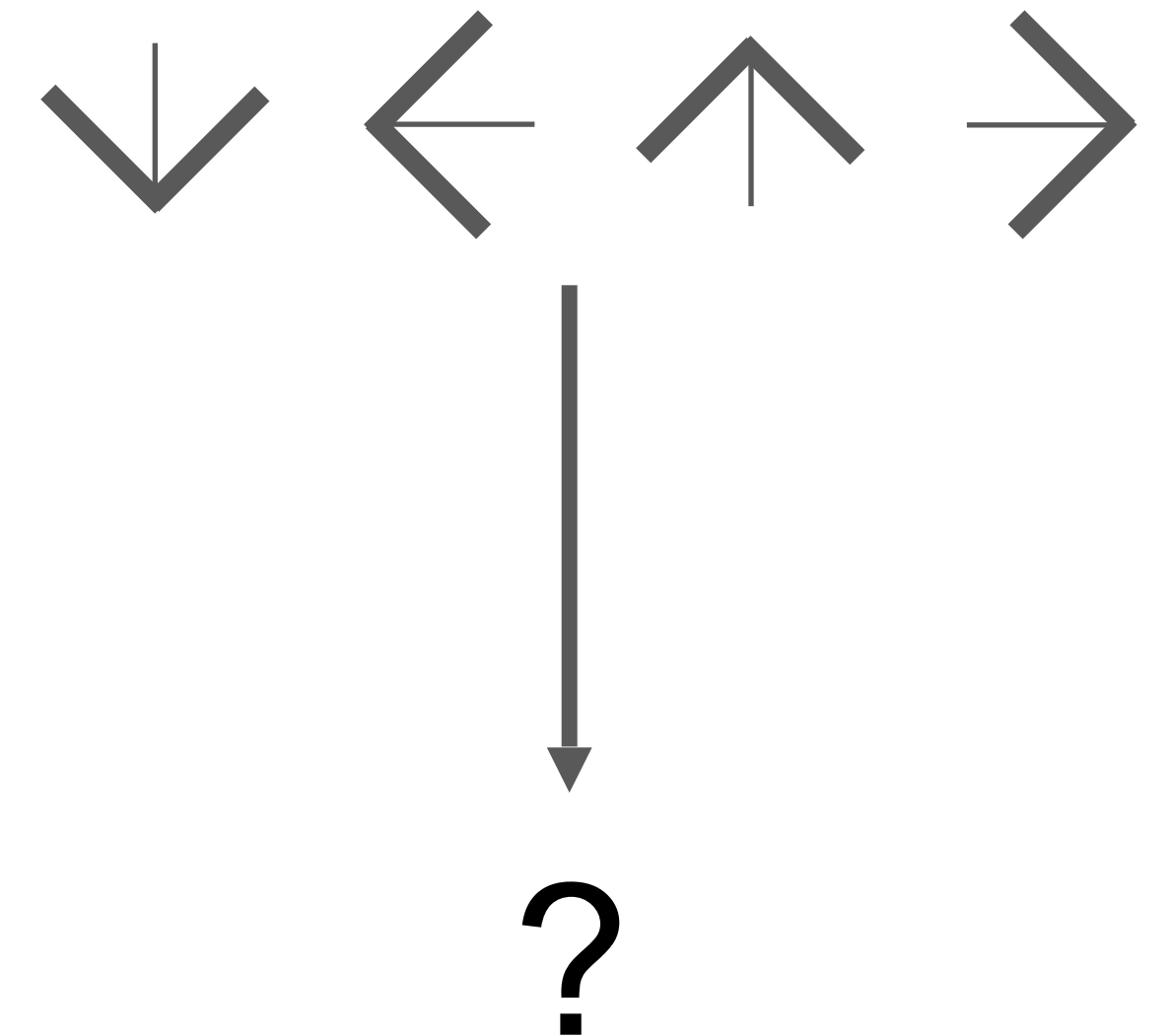**no obvious way** to hand-code the algorithm for recognizing a cat, or other classes.

# Attempts have been made



Find edges

Find corners

?

John Canny, "A Computational Approach to Edge Detection", IEEE TPAMI 1986

# Machine Learning: Data Driven Approach

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

**Example training set**

```
def train(train_images, train_labels):
    # build a model for images -> labels...
    return model


def predict(model, test_images):
    # predict test_labels using the model...
    return test_labels
```

Subhransu Maji — UMass Amherst, Spring 25

# Today

Examples of machine learning models
- Nearest neighbor classifiers
- Linear classifiers
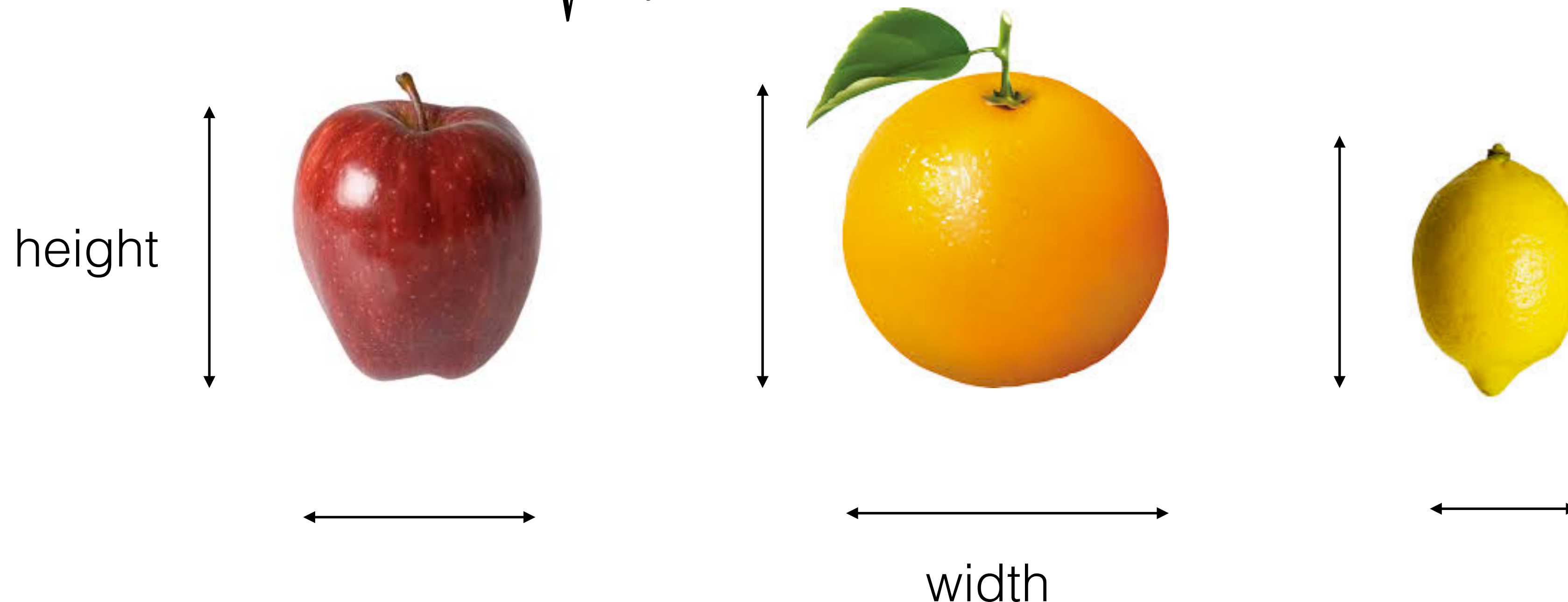
# Nearest Neighbor Classifier

# Nearest neighbor classifier

Training data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_n, y_n)$

Fruit data:

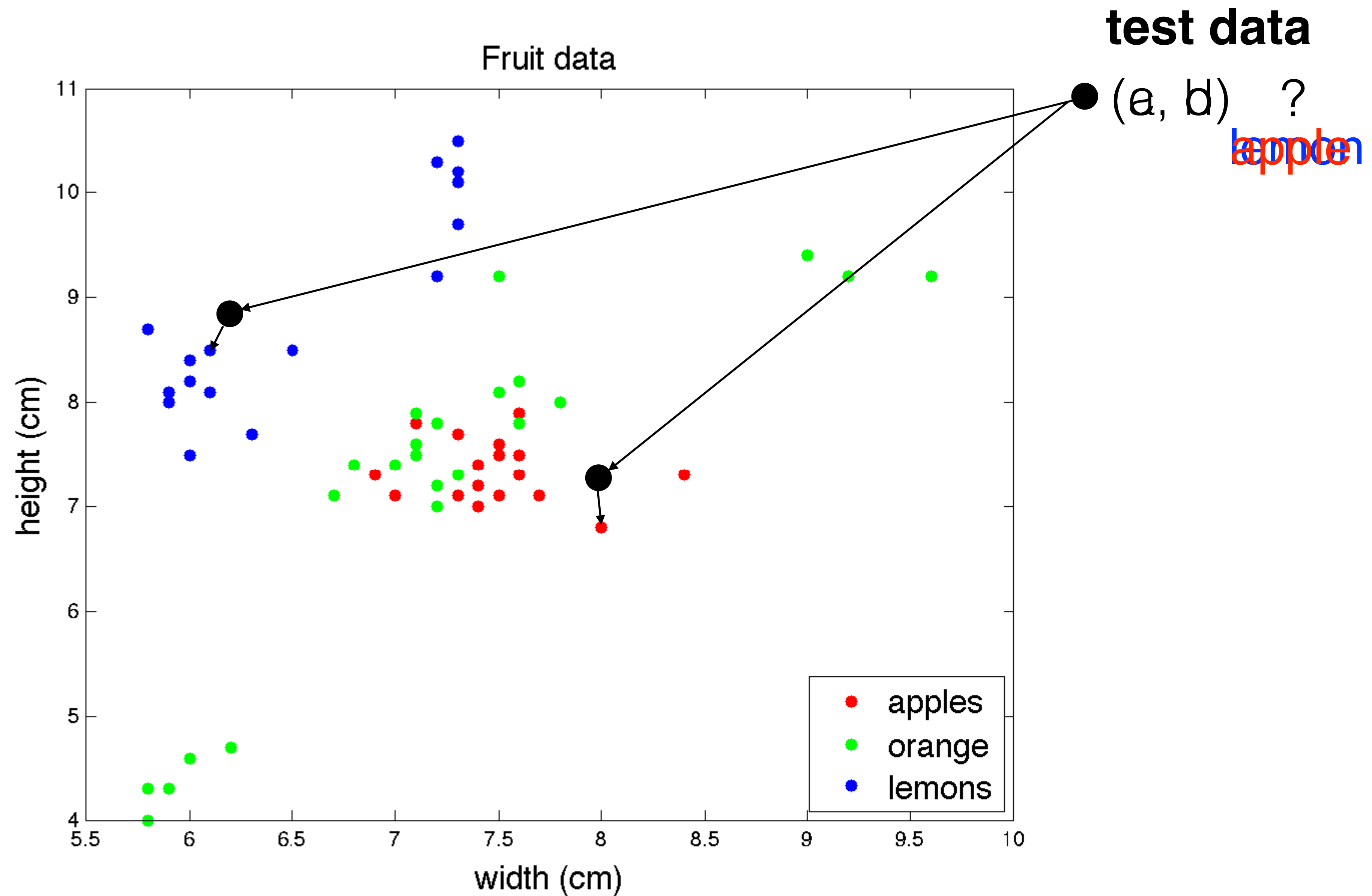- label: {apples, oranges, lemons}
- attributes: {width, height}

Euclidean distance $\quad d(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{\sum_i (\mathbf{x}_{1,i} - \mathbf{x}_{2,i})^2}$



height

width

# Nearest neighbor classifier

**test data**

(a, b)    ?

lemon apple



Fruit data

height (cm) vs width (cm)

Legend:
- apples (red)
- orange (green)
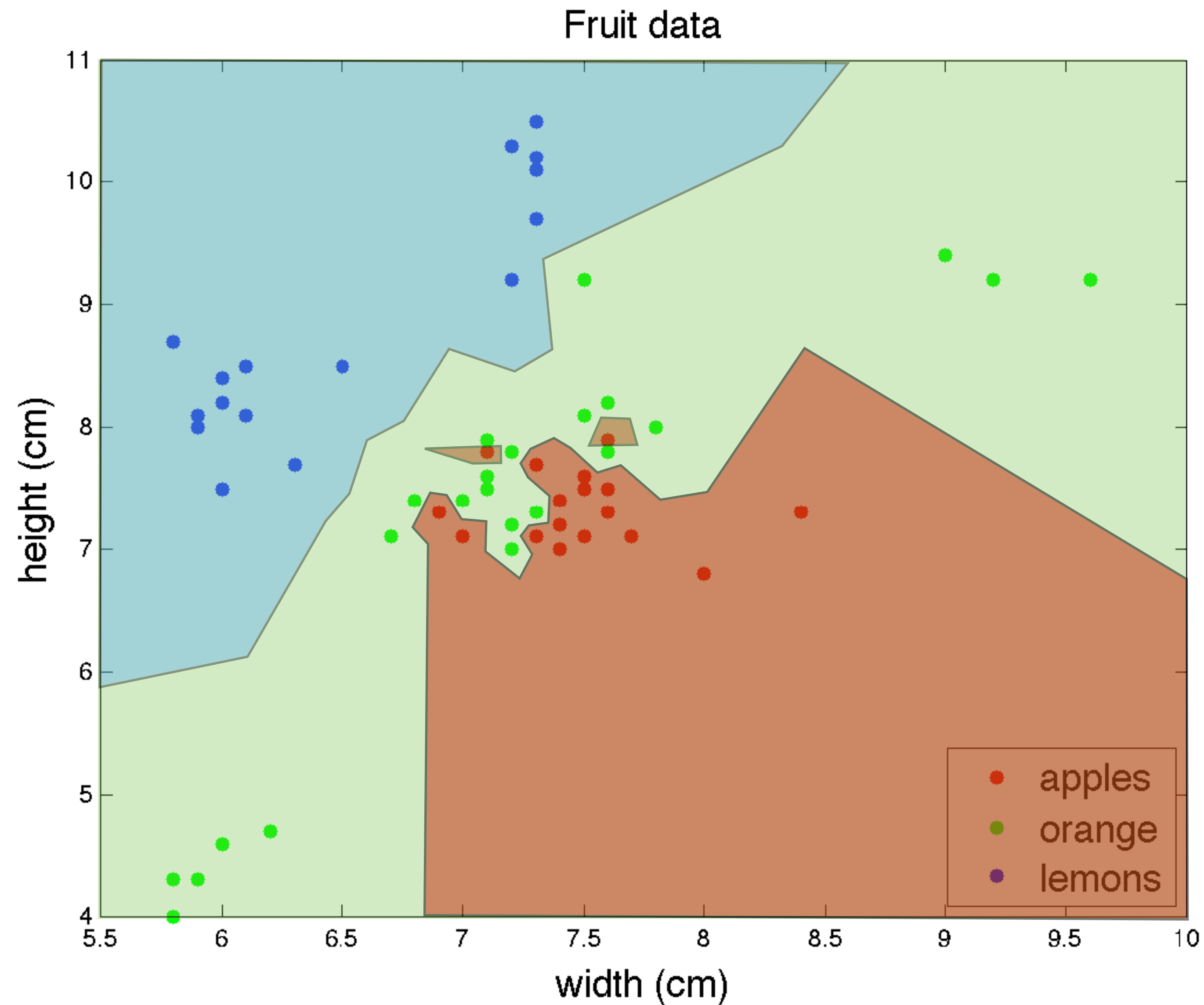- lemons (blue)

# Decision boundaries: 1NN



Fruit data

# k-Nearest neighbor classifier

Take majority vote among the k nearest neighbors



Fruit data

outlier

What is the effect of k?

- apples
- orange
- lemons

# k-Nearest Neighbor
## find the k nearest images, have them vote on the label

the data

NN classifier

5-NN classifier



http://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm

the data     NN classifier     5-NN classifier

**Q:** what is the accuracy of the nearest neighbor classifier on the training data, when using the Euclidean distance?

# Example dataset: **CIFAR-10**

**10** labels
**50,000** training images
**10,000** test images

For every test image (first column), examples of nearest neighbors in rows

# Nearest Neighbor classifier

```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```
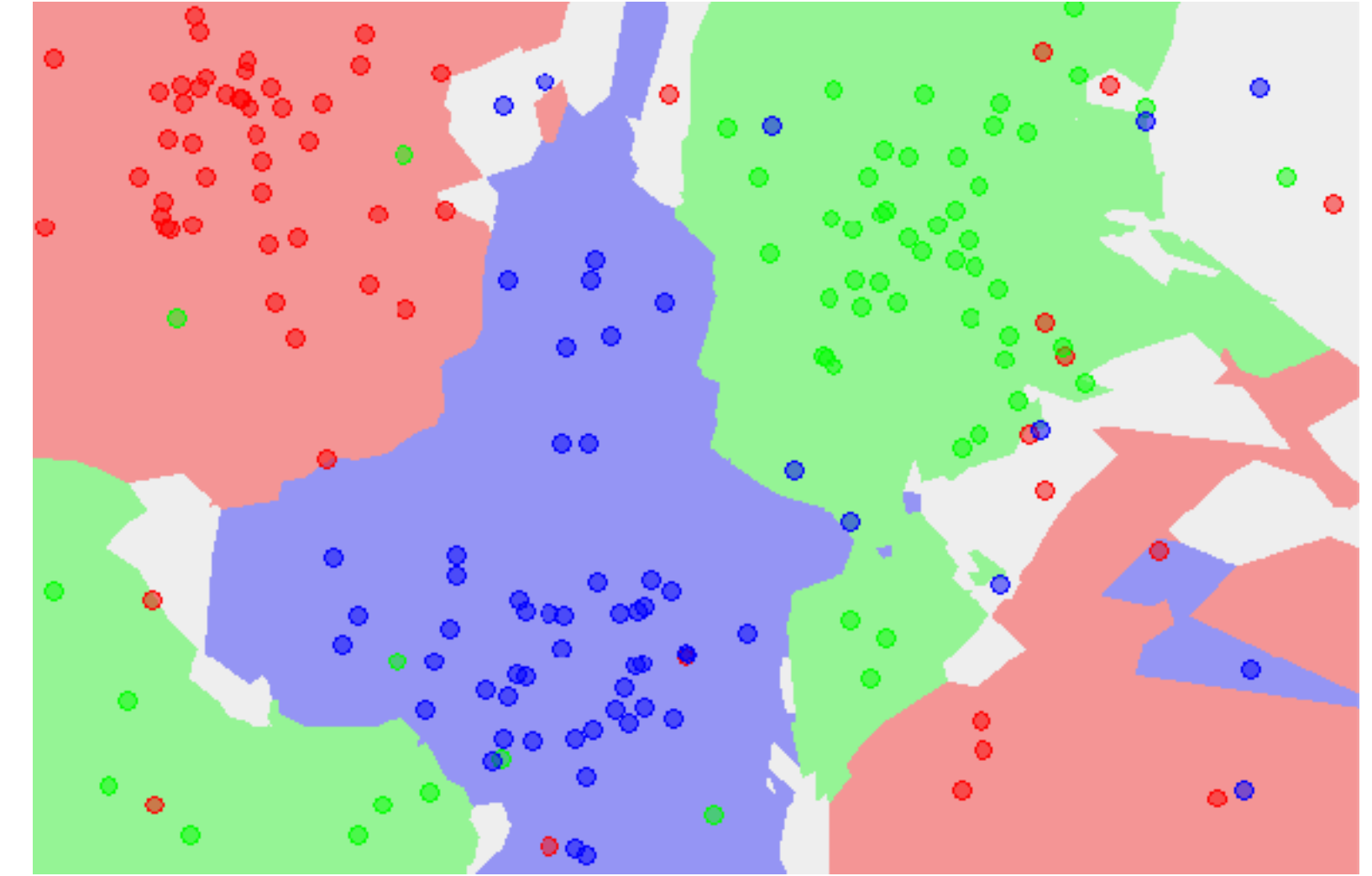
# Nearest Neighbor classifier

```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```

remember the training data

# Nearest Neighbor classifier

```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```

for every test image:
- find nearest train image with L1 distance
- predict the label of nearest training image

the data | NN classifier | 5-NN classifier

**Q:** Suppose you have N training examples. How long does it take to make a prediction with a nearest neighbor classifier on one test example?

What is the best **distance** to use?
What is the best value of **k** to use?

i.e. how do we set the **hyperparameters**?

What is the best **distance** to use?
What is the best value of **k** to use?

i.e. how do we set the **hyperparameters**?

Very problem-dependent.
Must try them all out and see what works best.

# Trying out what hyperparameters work best on test set.

Trying out what hyperparameters work best on test set:
Very bad idea. The test set is a proxy for the generalization performance!
Use only **VERY SPARINGLY,** at the end.



train data | test data

Validation data
use to tune hyperparameters

**Cross-validation**
cycle through the choice of which fold is the validation fold, average results.

Cross-validation on k

Example of
5-fold cross-validation
for the value of **k.**

Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that k ~= 7 works best
for this data)

# k-Nearest Neighbor on *raw* images is **never used.**

- terrible performance at test time
- distance metrics on level of whole images can be very unintuitive



| original | shifted | messed up | darkened |

(all 3 images have same L2 distance to the one on the left)

# So far …

Nearest neighbor classifier

- All features are equally good
- No training required!
- Slow at test time

Linear classifiers (next)

- Use all features, but some more than others
- Training required
- Fast at test time!

# Linear Classification

airplane

automobile

bird

cat

deer

dog

frog

horse

ship

truck

Example dataset: **CIFAR-10**
**10** labels
**50,000** training images
  each image is **32x32x3**
**10,000** test images.

# Parametric approach



image    parameters

$$f(\mathbf{x}, \mathbf{W})$$

**[32x32x3]**
array of numbers 0...1
(3072 numbers total)

**10** numbers, indicating class scores

# Parametric approach: **Linear classifier**

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



**[32x32x3]**
array of numbers 0...1

**10** numbers, indicating class scores

# Parametric approach: **Linear classifier**

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

**3072x1**

**10x1**   **10x3072**

[32x32x3]
array of numbers 0...1

→ **10** numbers, indicating class scores

parameters, or "weights"

# Parametric approach: **Linear classifier**

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

10x1    10x3072    3072x1    10x1



**[32x32x3]**
array of numbers 0...1

**10** numbers, indicating class scores

parameters, or "weights"

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

Flatten tensors into a vector



Input image

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

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



stretch pixels into single column

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

$W$

input image

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

$x_i$

**+**

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

$b$

$\longrightarrow$

| | |
|---|---|
| -96.8 | cat score |
| 437.9 | dog score |
| 61.95 | ship score |

$f(x_i; W, b)$

# Interpreting a Linear Classifier



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

Q: what does the linear classifier do, in English?

# Interpreting a Linear Classifier



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

**[32x32x3]**
array of numbers 0...1
(3072 numbers total)

# Interpreting a Linear Classifier



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

Example trained weights of a linear classifier trained on CIFAR-10:

# Interpreting a Linear Classifier



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

Q2: what would be a very hard set of classes for a linear classifier to distinguish?

# Hard cases for a linear classifier

**Class 1**:
First and third quadrants

**Class 2**:
Second and fourth quadrants

**Class 1**:
1 <= L2 norm <= 2

**Class 2**:
Everything else

**Class 1**:
Three modes

**Class 2**:
Everything else

# So far: We defined a (linear) **score function**: $f(x_i, W, b) = Wx_i + b$

really *affine*



Example class scores for 3 images, with a random W:

| | cat | car | frog |
|---|---|---|---|
| airplane | −3.45 | −0.51 | 3.42 |
| automobile | −8.87 | **6.04** | 4.64 |
| bird | 0.09 | 5.31 | 2.65 |
| cat | **2.9** | −4.22 | 5.1 |
| deer | 4.48 | −4.19 | 2.64 |
| dog | 8.02 | 3.58 | 5.55 |
| frog | 3.78 | 4.49 | **−4.34** |
| horse | 1.06 | −4.37 | −1.5 |
| ship | −0.36 | −2.09 | −4.79 |
| truck | −0.72 | −2.93 | 6.14 |

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

# Coming up:
- ## Loss function (quantifying what it means to have a "good" W)

- ## Optimization (start with random W and find a W that minimizes the loss)

- ## Neural nets! (tweak the functional form of f)

# Summary so far ...   Linear classifier



**image   parameters**

$$f(\mathbf{x}, \mathbf{W})$$

[32x32x3]
array of numbers 0...1
(3072 numbers total)

**10** numbers, indicating class scores

stretch pixels into single column

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

$W$

input image

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

$x_i$

**+**

| 1.1 |
| 3.2 |
| -1.2 |

$b$

→

| -96.8 | cat score |
| 437.9 | dog score |
| 61.95 | ship score |

$f(x_i; W, b)$

airplane classifier
car classifier
deer classifier

plane  car  bird  cat  deer  dog  frog  horse  ship  truck

# Loss function/Optimization



| | cat | car | frog |
|---|---|---|---|
| airplane | -3.45 | -0.51 | 3.42 |
| automobile | -8.87 | **6.04** | 4.64 |
| bird | 0.09 | 5.31 | 2.65 |
| cat | **2.9** | -4.22 | 5.1 |
| deer | 4.48 | -4.19 | 2.64 |
| dog | 8.02 | 3.58 | 5.55 |
| frog | 3.78 | 4.49 | **-4.34** |
| horse | 1.06 | -4.37 | -1.5 |
| ship | -0.36 | -2.09 | -4.79 |
| truck | -0.72 | -2.93 | 6.14 |

TODO:

1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.

1. Come up with a way of efficiently finding the parameters that minimize the loss function. **(optimization)**

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** |

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** |

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

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

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

|  | 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 | | |

= 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    | 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      |        |

**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:  $s_i = f(x_i, W)$

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:

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

|       | 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 |

= max(0, 2.2 - (-3.1) + 1)
   +max(0, 2.5 - (-3.1) + 1)
= max(0, 6.3) + max(0, 6.6)
= 6.3 + 6.6
= 12.9

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 | 0 | 12.9 |

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

and the full training loss is the mean
over all examples in the training data:

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i$$

L = (2.9 + 0 + 12.9)/3
   = **5.3**

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 |

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

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

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 |

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

Q2: what if we used a mean instead of a sum here?

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 |

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

Q3: what if we used

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$$

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 |

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

Q4: what is the min/max possible loss?

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 |

**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: $s_i = f(x_i, W)$

the SVM loss has the form:

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

Q5: usually at initialization W are small numbers, so all s ~= 0. What is the loss?

# Example numpy code:

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

```python
def L_i_vectorized(x, y, W):
    scores = W.dot(x)
    margins = np.maximum(0, scores - scores[y] + 1)
    margins[y] = 0
    loss_i = np.sum(margins)
    return loss_i
```

# Coding tip: Keep track of dimensions:

```
N = X.shape[0]
D = X.shape[1]
C = W.shape[1]

scores=X.dot(W)              # (N,D)*(D,C)=(N,C)
```

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



cat     **3.2**

car     5.1

frog    -1.7

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



**scores = unnormalized log probabilities of the classes.**

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

cat     **3.2**

car     5.1

frog    -1.7

# Softmax Classifier (Multinomial Logistic Regression)

**scores = unnormalized log probabilities of the classes.**

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

where

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

cat **3.2**

car 5.1

frog -1.7

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



**scores = unnormalized log probabilities of the classes.**

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

where

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

<span style="color:red">Softmax function</span>

cat **3.2**

car 5.1

frog -1.7

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



**scores = unnormalized log probabilities of the classes.**

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ where $$s = f(x_i; W)$$

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

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

cat **3.2**

car 5.1

frog -1.7

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



cat **3.2**

car 5.1

frog -1.7

**scores = unnormalized log probabilities of the classes.**

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ where $s = f(x_i; W)$

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

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

in summary: $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$

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

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

cat

car

frog

| |
|---|
| **3.2** |
| 5.1 |
| -1.7 |

<span style="color:blue">unnormalized log probabilities</span>

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

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

unnormalized probabilities

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

unnormalized log probabilities

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

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



unnormalized probabilities

|  | | exp | | normalize | |
|---|---|---|---|---|---|
| cat | **3.2** | → | **24.5** | → | **0.13** |
| car | 5.1 | | 164.0 | | 0.87 |
| frog | -1.7 | | 0.18 | | 0.00 |

unnormalized log probabilities

probabilities
>0, sum to 1

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

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

unnormalized probabilities

|  | unnormalized log probabilities | exp | | normalize | probabilities | |
|---|---|---|---|---|---|---|
| cat | **3.2** | → | **24.5** | → | **0.13** | → L_i = -log(0.13) = **0.89** |
| car | 5.1 | | 164.0 | | 0.87 | |
| frog | -1.7 | | 0.18 | | 0.00 | |

unnormalized log probabilities

probabilities

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

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

unnormalized probabilities

|       | unnormalized log probabilities | exp | | normalize | probabilities | |
|-------|--------------------------------|-----|--|-----------|---------------|--|
| cat   | **3.2**                        | →   | **24.5**  | →   | **0.13** | → L_i = -log(0.13) = **0.89** |
| car   | 5.1                            |     | 164.0     |     | 0.87     |   |
| frog  | -1.7                           |     | 0.18      |     | 0.00     |   |

unnormalized log probabilities

probabilities

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



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

unnormalized probabilities

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

| | unnormalized log probabilities | | unnormalized probabilities | | probabilities | |
|---|---|---|---|---|---|---|
| cat | **3.2** | exp → | **24.5** | normalize → | **0.13** | → L_i = -log(0.13) = **0.89** |
| car | 5.1 | | 164.0 | | 0.87 | |
| frog | -1.7 | | 0.18 | | 0.00 | |

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

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

unnormalized probabilities

Q2: usually at initialization W are small numbers, so all s ~= 0. What is the loss?

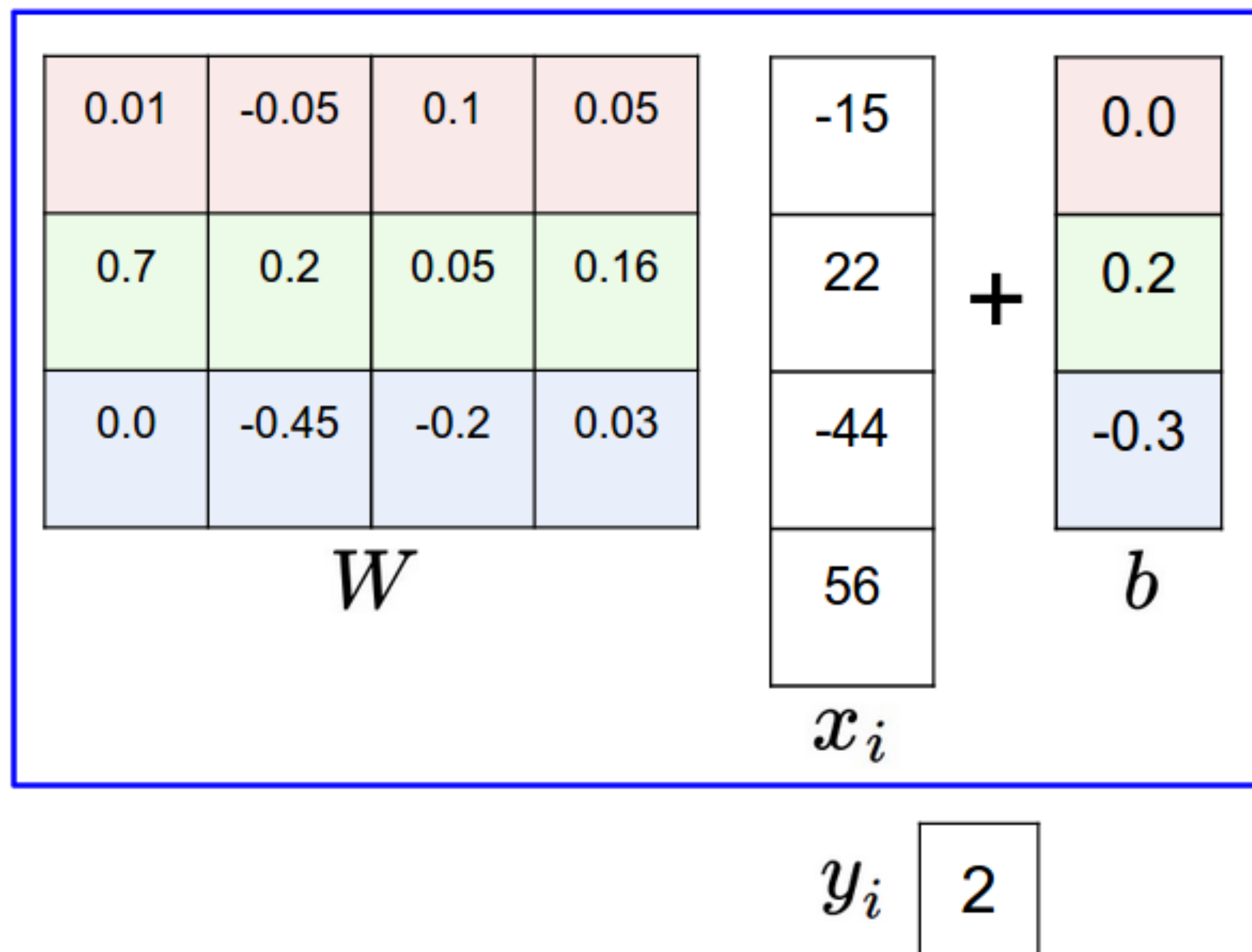|  | unnormalized log probabilities | | exp | unnormalized probabilities | | normalize | probabilities | |
|---|---|---|---|---|---|---|---|---|
| cat | **3.2** | | | **24.5** | | | **0.13** | L_i = -log(0.13) = **0.89** |
| car | 5.1 | | | 164.0 | | | 0.87 | |
| frog | -1.7 | | | 0.18 | | | 0.00 | |

# Softmax vs. SVM

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

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

# Softmax vs. SVM
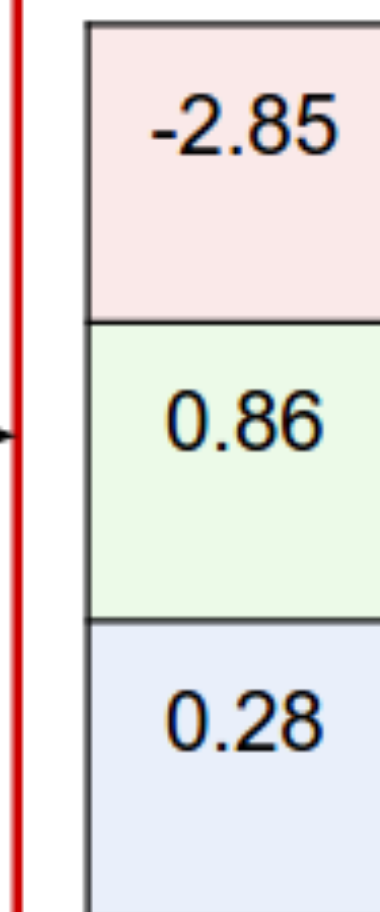
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

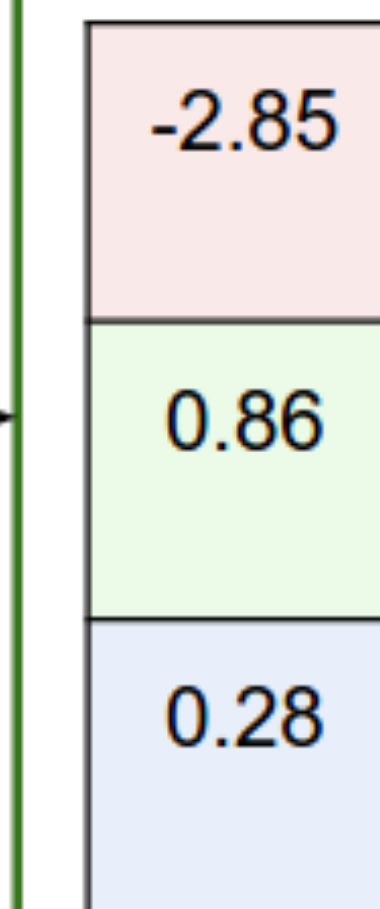$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

assume scores:
[10, -2, 3]
[10, 9, 9]
[10, -100, -100]
and  $\boxed{y_i = 0}$

Q: Suppose I take a datapoint and I jiggle a bit (changing its score slightly). What happens to the loss in both cases?

# Coming up:

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

# - Regularization
# - Optimization

# Regularization

There is a "bug" with the loss:

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

$$L = \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

E.g. Suppose that we found a W such that L = 0.
Is this W unique?

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

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

|  | 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 |

**Before:**
= 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

**With W twice as large:**
= max(0, 2.6 - 9.8 + 1)
    +max(0, 4.0 - 9.8 + 1)
= max(0, -6.2) + max(0, -4.8)
= 0 + 0
= 0

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



An example:

What is the loss? (POLL)

cat       1.3

car       **2.5**

frog      2.0

Loss:

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



An example:
What is the loss?

| | |
|---|---|
| cat | 1.3 |
| car | **2.5** |
| frog | 2.0 |
| Loss: | 0.5 |

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



An example:

　　　What is the loss?

How could we change W to eliminate the loss?  (POLL)

| | |
|---|---|
| cat | 1.3 |
| car | **2.5** |
| frog | 2.0 |
| Loss: | 0.5 |

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



An example:

What is the loss?

How could we change W to eliminate the loss? (POLL)

Multiply W (and b) by 2!

| | | |
|---|---|---|
| cat | 1.3 | 2.6 |
| car | **2.5** | **5.0** |
| frog | 2.0 | 4.0 |
| Loss: | 0.5 | 0 |

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



An example:

What is the loss?

How could we change W to eliminate the loss? (POLL)

Multiply W (and b) by 2!

Wait a minute! Have we done anything useful???

|       |       |       |
|-------|-------|-------|
| cat   | 1.3   | 2.6   |
| car   | **2.5** | **5.0** |
| frog  | 2.0   | 4.0   |
| Loss: | 0.5   | 0     |

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

An example:

What is the loss?

How could we change W to eliminate the loss? (POLL)

Multiply W (and b) by 2!

Wait a minute! Have we done anything useful???

No! Any example that used to be wrong is still wrong (on the wrong side of the boundary). Any example that is right is still right (on the correct side of the boundary).

| | | |
|---|---|---|
| cat | 1.3 | 2.6 |
| car | **2.5** | **5.0** |
| frog | 2.0 | 4.0 |
| Loss: | 0.5 | 0 |

# Regularization

$\lambda$ = regularization strength (hyperparameter)

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

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from having too much flexibility.

**Simple examples**

L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

**More complex**:

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc

# Regularization

$\lambda$ = regularization strength (hyperparameter)

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

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from having too much flexibility.

Why regularize?
- Express preferences over weights
- Make the model *simple* so it works on test data
- Improve optimization by adding curvature

# Regularization: Expressing Preferences

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

$$w_1^T x = w_2^T x = 1$$

L2 Regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

# Regularization: Expressing Preferences

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = \boxed{[0.25, 0.25, 0.25, 0.25]}$$

$$w_1^T x = w_2^T x = 1$$

L2 Regularization

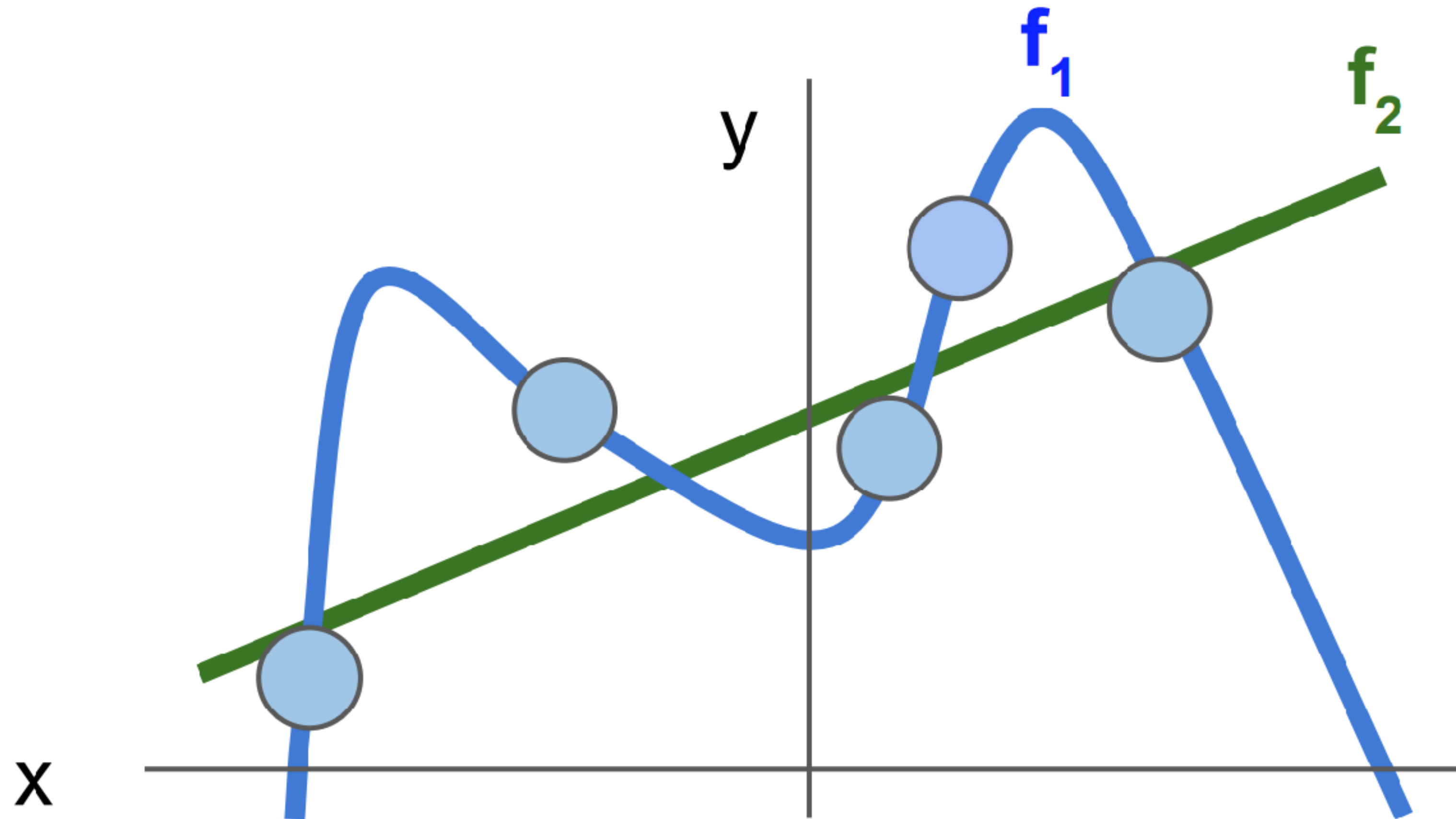$$R(W) = \sum_k \sum_l W_{k,l}^2$$

L2 regularization likes to "spread out" the weights

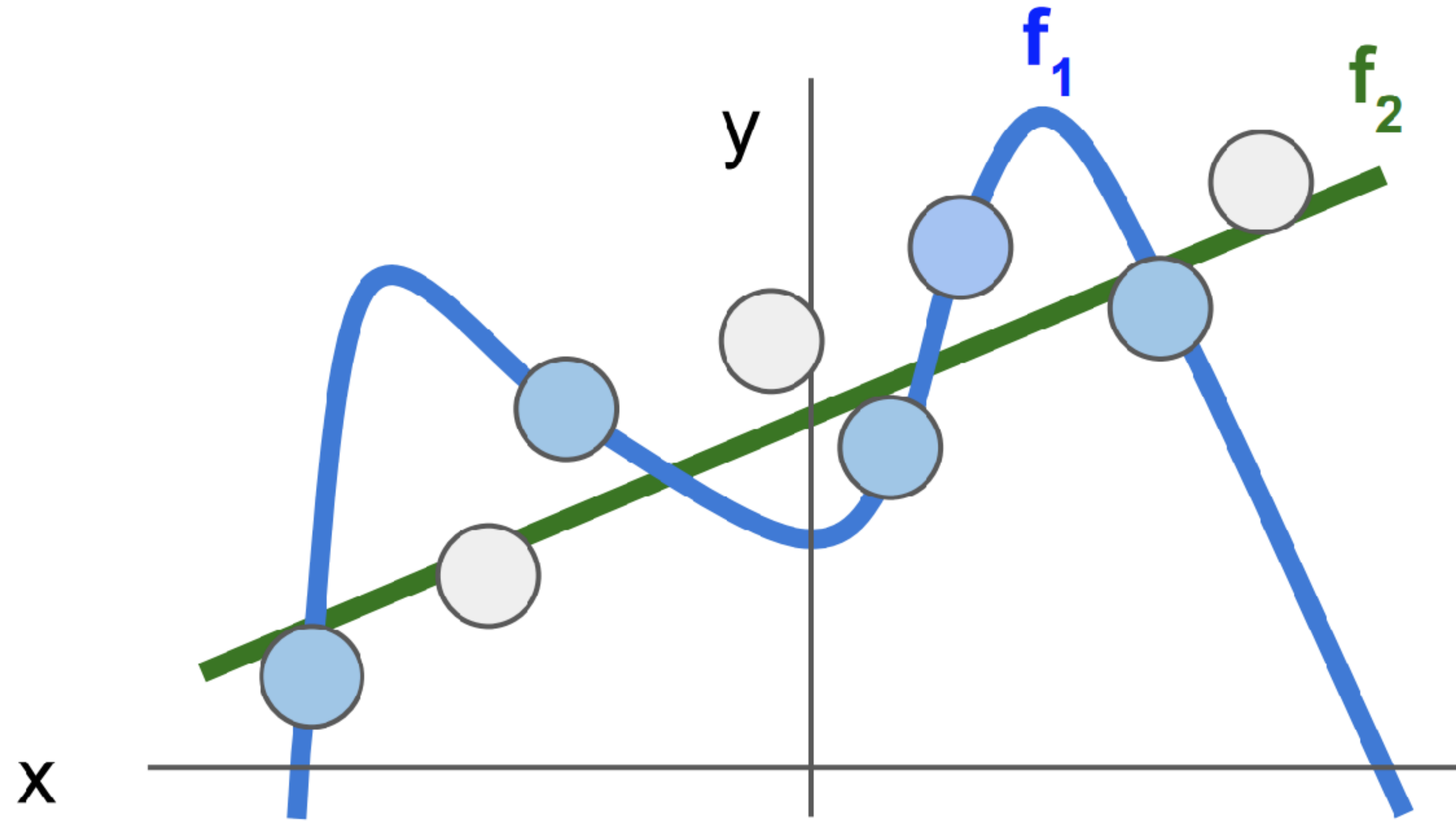# Regularization: Prefer Simpler Models

# Regularization: Prefer Simpler Models

# Regularization: Prefer Simpler Models



Regularization pushes against fitting the data with too much flexibility. If you are going to use a complex function to fit the data, you should be doing based on a lot of data!
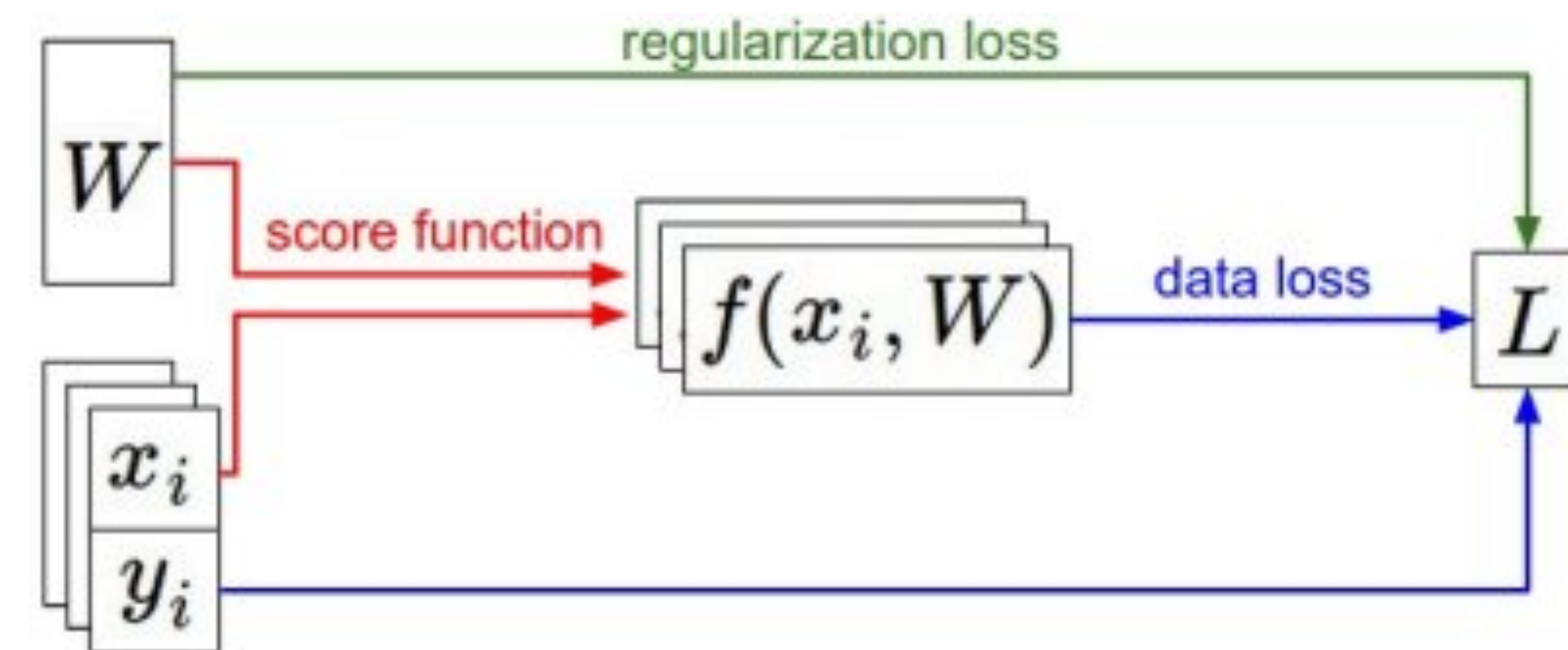
# Optimization

# Recap

- We have some dataset of (x,y)
- We have a **score function:**
- We have a **loss function:**

$$s = f(x; W) \overset{\text{e.g.}}{=} Wx$$

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \text{ Softmax}$$

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

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + R(W) \quad \text{Full loss}$$

# Strategy #1: A first very bad idea solution: **Random search**

```python
# assume X_train is the data where each column is an example (e.g. 3073 x 50,000)
# assume Y_train are the labels (e.g. 1D array of 50,000)
# assume the function L evaluates the loss function

bestloss = float("inf") # Python assigns the highest possible float value
for num in xrange(1000):
  W = np.random.randn(10, 3073) * 0.0001 # generate random parameters
  loss = L(X_train, Y_train, W) # get the loss over the entire training set
  if loss < bestloss: # keep track of the best solution
    bestloss = loss
    bestW = W
  print 'in attempt %d the loss was %f, best %f' % (num, loss, bestloss)

# prints:
# in attempt 0 the loss was 9.401632, best 9.401632
# in attempt 1 the loss was 8.959668, best 8.959668
# in attempt 2 the loss was 9.044034, best 8.959668
# in attempt 3 the loss was 9.278948, best 8.959668
# in attempt 4 the loss was 8.857370, best 8.857370
# in attempt 5 the loss was 8.943151, best 8.857370
# in attempt 6 the loss was 8.605604, best 8.605604
# ... (trunctated: continues for 1000 lines)
```

# Let's see how well this works on the test set...

```python
# Assume X_test is [3073 x 10000], Y_test [10000 x 1]
scores = Wbest.dot(Xte_cols) # 10 x 10000, the class scores for all test examples
# find the index with max score in each column (the predicted class)
Yte_predict = np.argmax(scores, axis = 0)
# and calculate accuracy (fraction of predictions that are correct)
np.mean(Yte_predict == Yte)
# returns 0.1555
```
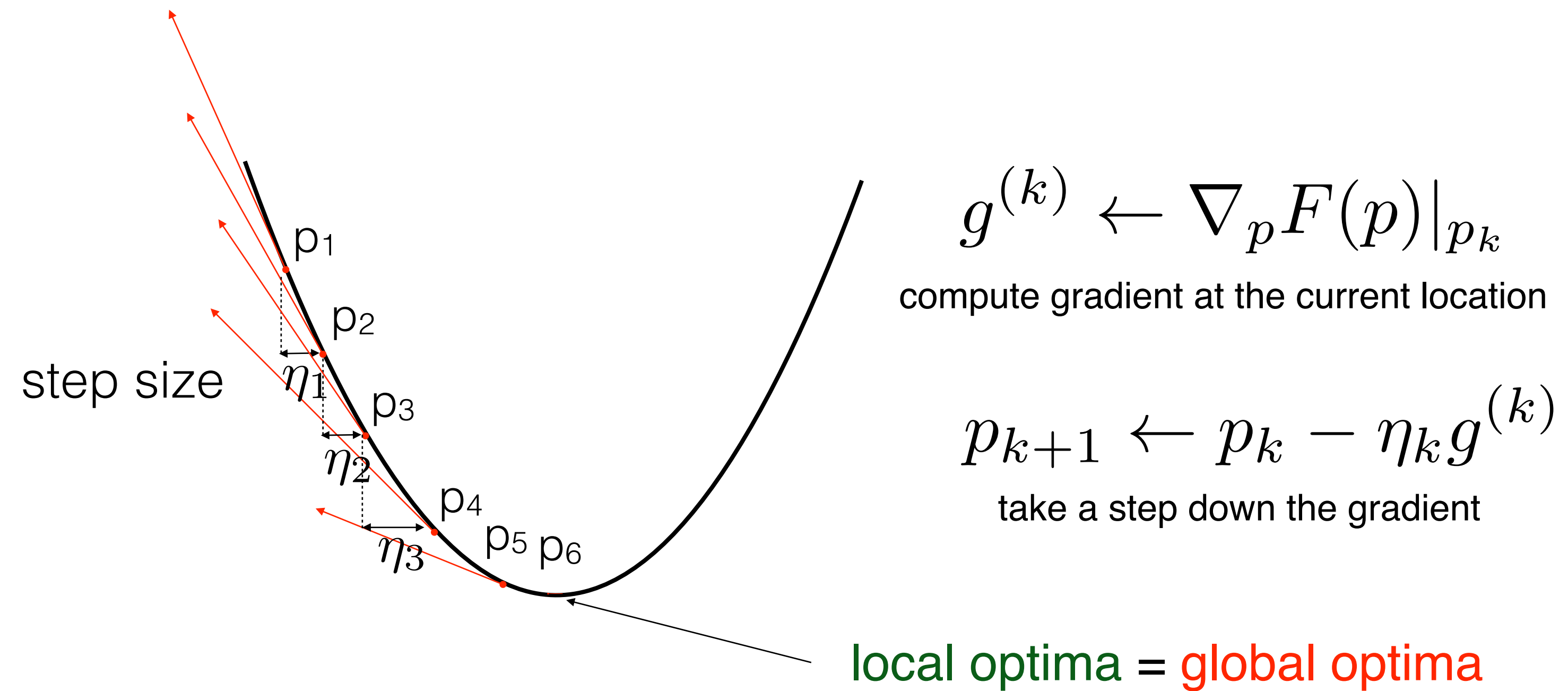
15.5% accuracy! not bad!
(SOTA is ~95%)

How often will a random search succeed?

# Strategy #2: **Follow the slope**



$$g^{(k)} \leftarrow \nabla_p F(p)\big|_{p_k}$$

compute gradient at the current location

$$p_{k+1} \leftarrow p_k - \eta_k g^{(k)}$$

take a step down the gradient

step size

p1
p2
$\eta_1$
p3
$\eta_2$
p4
$\eta_3$
p5  p6

local optima = global optima

# Strategy #2: **Follow the slope**

In 1-dimension, the derivative of a function:

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

In multiple dimensions, the **gradient** is the vector of (partial derivatives).

# Numerical evaluation of the gradient...

**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:** | **W + h** (second dim): | **gradient dW:** |
|---|---|---|
| [0.34, | [0.34, | [-2.5, |
| -1.11, | -1.11 + **0.0001**, | ?, |
| 0.78, | 0.78, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33,…] | 0.33,…] | ?,…] |
| **loss 1.25347** | **loss 1.25353** | |

**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:** | **W + h** (third dim): |
|---|---|
| [0.34, | [0.34, |
| -1.11, | -1.11, |
| 0.78, | 0.78 + **0.0001**, |
| 0.12, | 0.12, |
| 0.55, | 0.55, |
| 2.81, | 2.81, |
| -3.1, | -3.1, |
| -1.5, | -1.5, |
| 0.33,…] | 0.33,…] |
| **loss 1.25347** | **loss 1.25347** |

**gradient dW:**
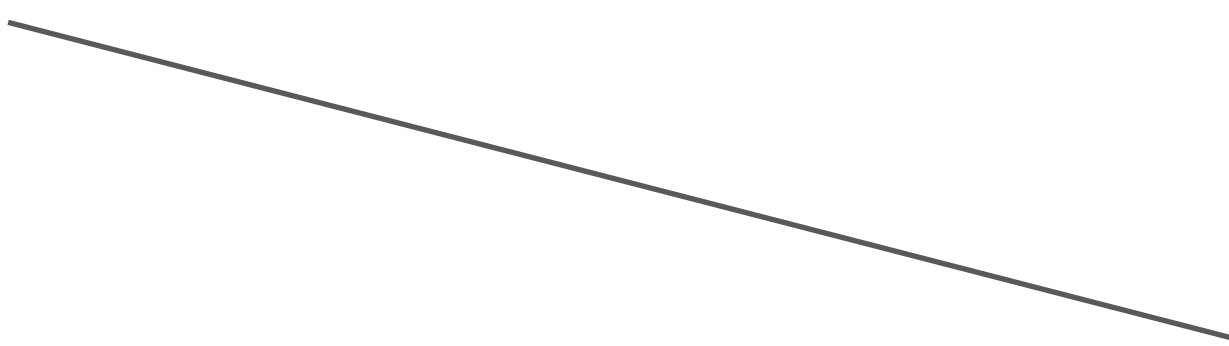
[-2.5,
0.6,
**0**,
?,
?

(1.25347 - 1.25347)/0.0001
= 0

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

?,…]

**current W:**

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

dW = ...
(some function of
data and W)

**gradient dW:**

[-2.5,
0.6,
0,
0.2,
0.7,
-0.5,
1.1,
1.3,
-2.1,…]

# Evaluating the gradient numerically

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

```python
def eval_numerical_gradient(f, x):
    """
    a naive implementation of numerical gradient of f at x
    - f should be a function that takes a single argument
    - x is the point (numpy array) to evaluate the gradient at
    """

    fx = f(x) # evaluate function value at original point
    grad = np.zeros(x.shape)
    h = 0.00001

    # iterate over all indexes in x
    it = np.nditer(x, flags=['multi_index'], op_flags=['readwrite'])
    while not it.finished:

        # evaluate function at x+h
        ix = it.multi_index
        old_value = x[ix]
        x[ix] = old_value + h # increment by h
        fxh = f(x) # evalute f(x + h)
        x[ix] = old_value # restore to previous value (very important!)

        # compute the partial derivative
        grad[ix] = (fxh - fx) / h # the slope
        it.iternext() # step to next dimension

    return grad
```

# Evaluating the gradient numerically

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

- approximate
- very slow to evaluate

```python
def eval_numerical_gradient(f, x):
  """
  a naive implementation of numerical gradient of f at x
  - f should be a function that takes a single argument
  - x is the point (numpy array) to evaluate the gradient at
  """

  fx = f(x) # evaluate function value at original point
  grad = np.zeros(x.shape)
  h = 0.00001

  # iterate over all indexes in x
  it = np.nditer(x, flags=['multi_index'], op_flags=['readwrite'])
  while not it.finished:

    # evaluate function at x+h
    ix = it.multi_index
    old_value = x[ix]
    x[ix] = old_value + h # increment by h
    fxh = f(x) # evalute f(x + h)
    x[ix] = old_value # restore to previous value (very important!)

    # compute the partial derivative
    grad[ix] = (fxh - fx) / h # the slope
    it.iternext() # step to next dimension

  return grad
```

# This is silly. The loss is just a function of W:

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \sum_k W_k^2$$

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

$$s = f(x; W) = Wx$$

want $\nabla_W L$ &larr;⎯⎯⎯⎯⎯⎯ "The gradient of the loss L with respect to the parameters W"

# This is silly. The loss is just a function of W:

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \sum_k W_k^2$$

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

$$s = f(x; W) = Wx$$

want $\nabla_W L$

# During a pandemic, Isaac Newton had to work from home, too. He used the time wisely.



A later portrait of Sir Isaac Newton by Samuel Freeman. (British Library/National Endowment for the Humanities)

By **Gillian Brockell**

March 12, 2020 at 2:18 p.m. EDT

Isaac Newton was in his early 20s when the Great Plague of London hit. He wasn't a "Sir" yet, didn't

1. Developed calculus
2. Fundamentals of optics
3. Theory of gravity

   ...not too shabby!

# This is silly. The loss is just a function of W:

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \sum_k W_k^2$$

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

$$s = f(x; W) = Wx$$

$$\nabla_W L \ = \ ...$$

# In summary:

- Numerical gradient: approximate, slow, easy to write

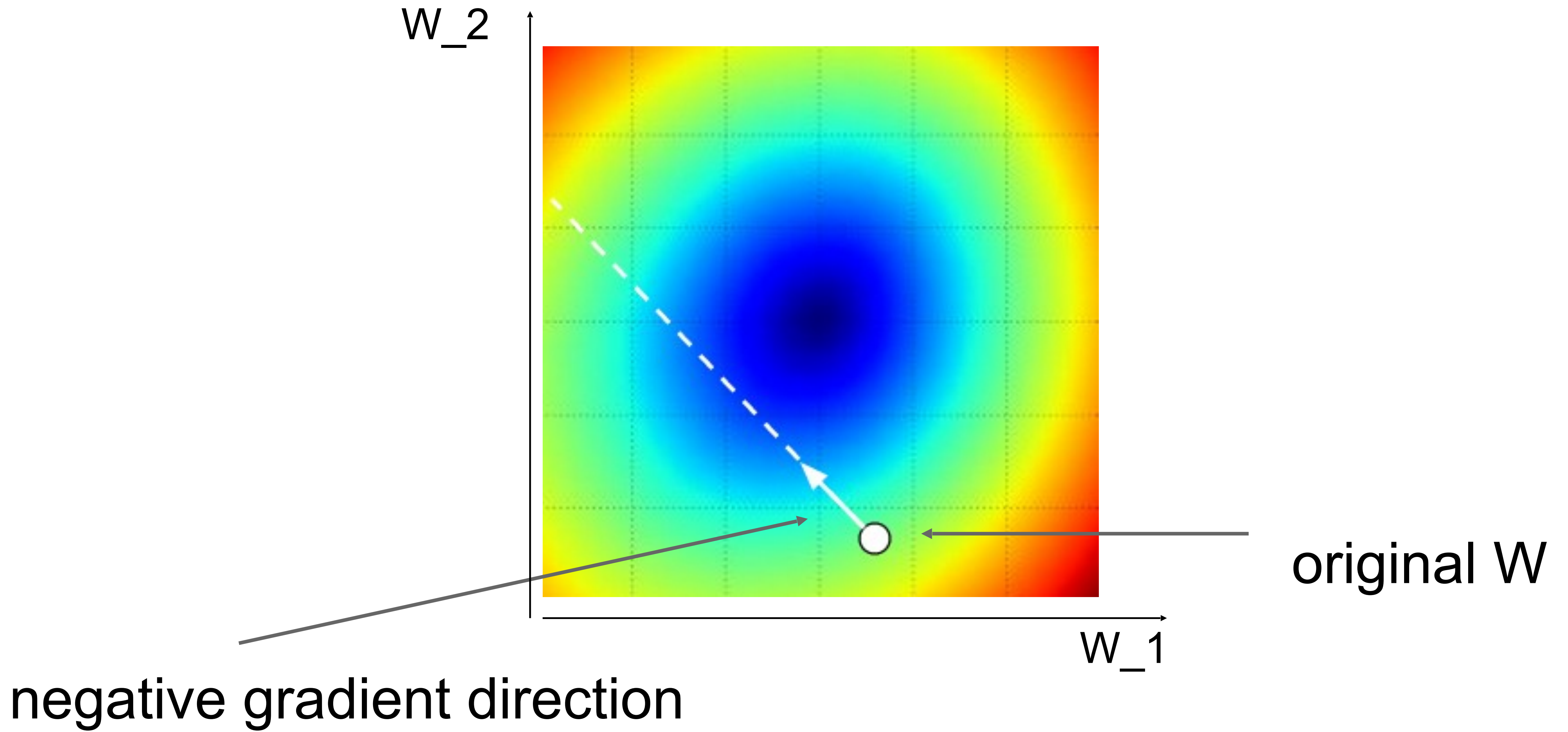- Analytic gradient: exact, fast, error-prone

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check.**

# Gradient Descent

```
# Vanilla Gradient Descent

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

W_2

W_1

original W

negative gradient direction

# Mini-batch Gradient Descent

- only use a small portion of the training set to compute the gradient.

```
# 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
```

Common mini-batch sizes are 32/64/128 examples
e.g. Krizhevsky ILSVRC ConvNet used 256 examples

# Mini-batch Gradient Descent

- only use a small portion of the training set to compute the gradient. Why?
    - Goal is to estimate the gradient
    - Trade-off between accuracy and computation
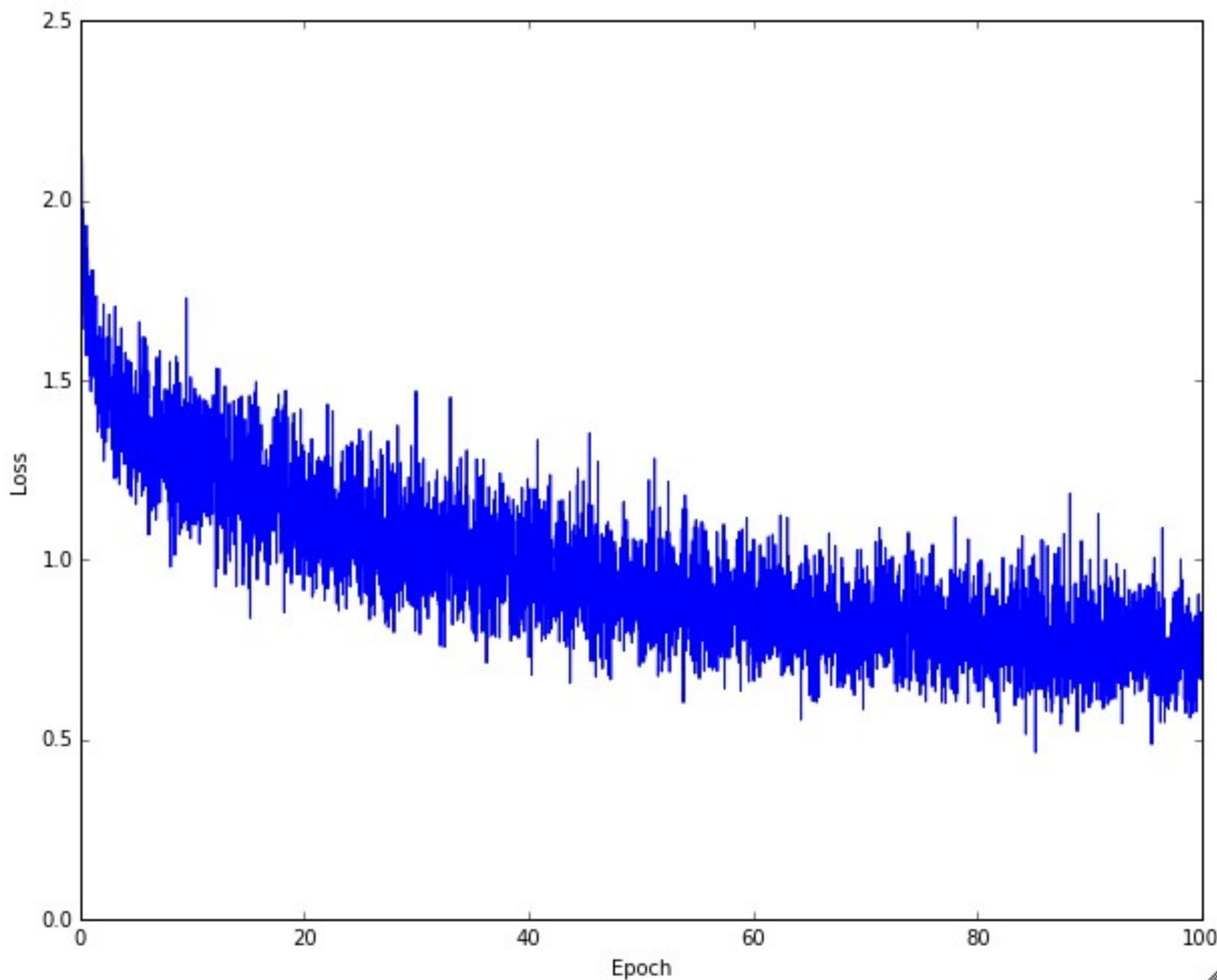    - No point in doing more computation if it won't change the updates

# Mini-batch Gradient Descent

- only use a small portion of the training set to compute the gradient.

```python
# 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
```

Common mini-batch sizes are 32/64/128 examples
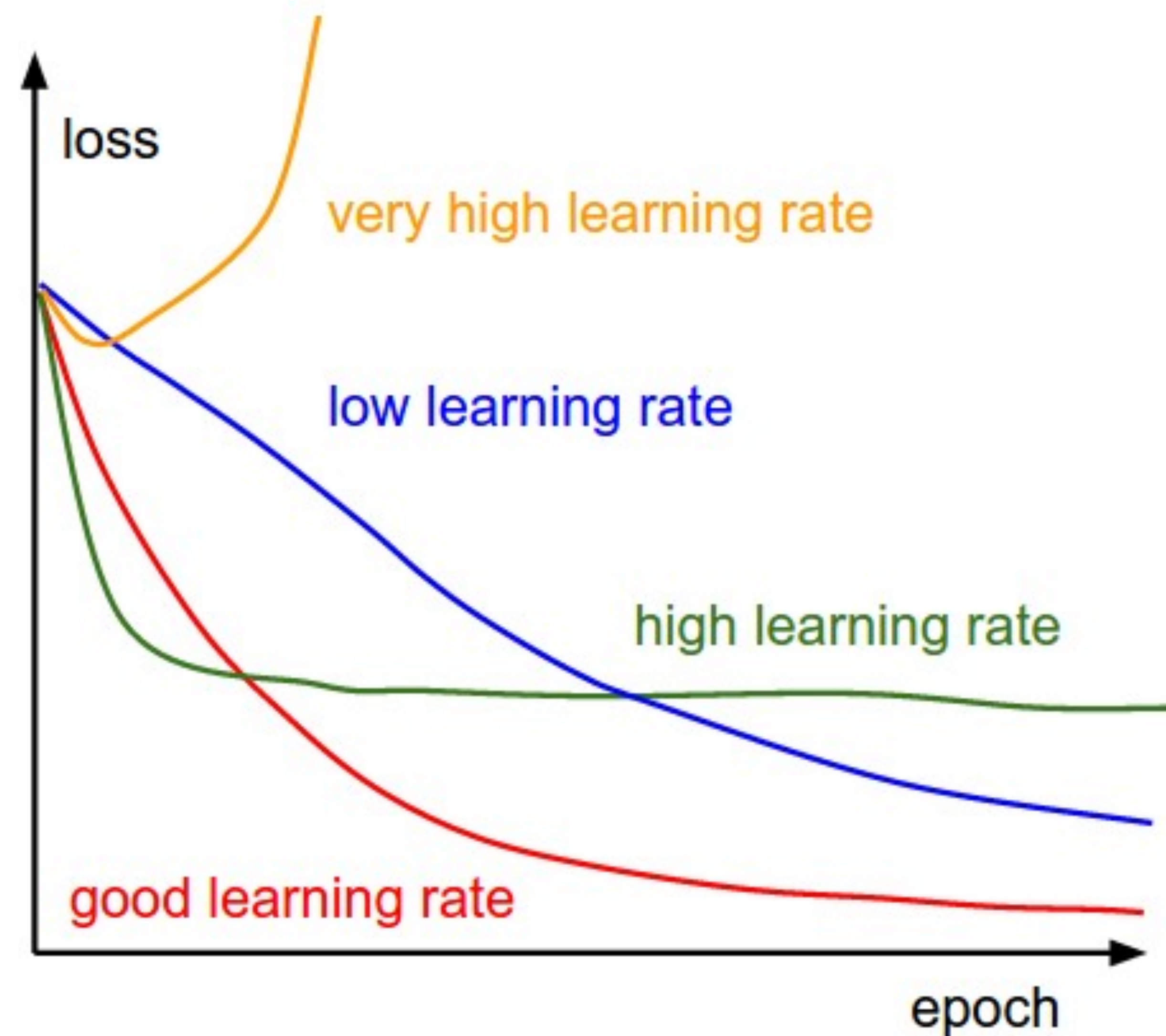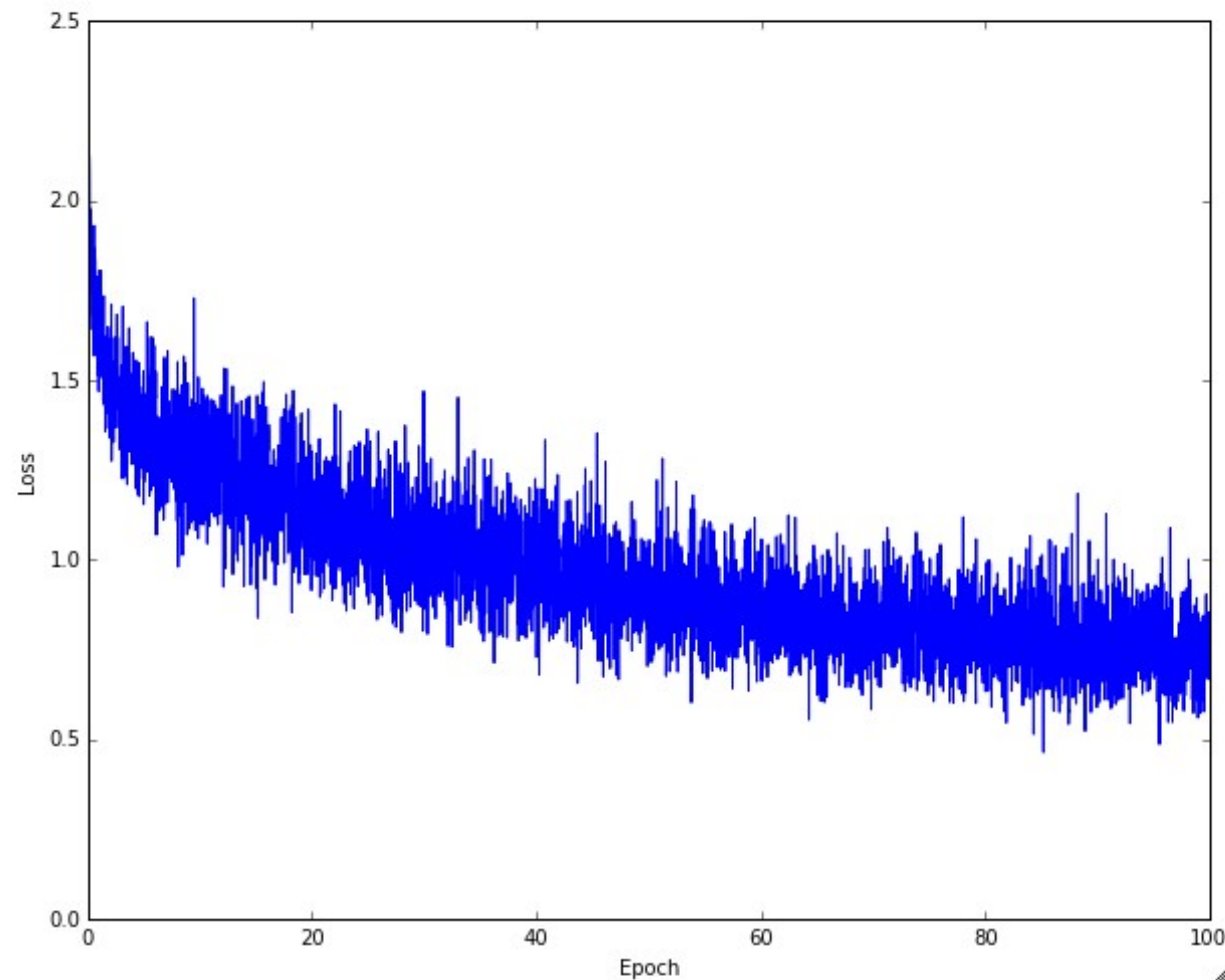e.g. Krizhevsky ILSVRC ConvNet used 256 examples

Example of optimization progress while training a neural network.

(Loss over mini-batches goes down over time.)

# The effects of step size (or "learning rate")





very high learning rate

loss

low learning rate

high learning rate

good learning rate

epoch

# Mini-batch Gradient Descent

- only use a small portion of the training set to compute the gradient.

```
# 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
```
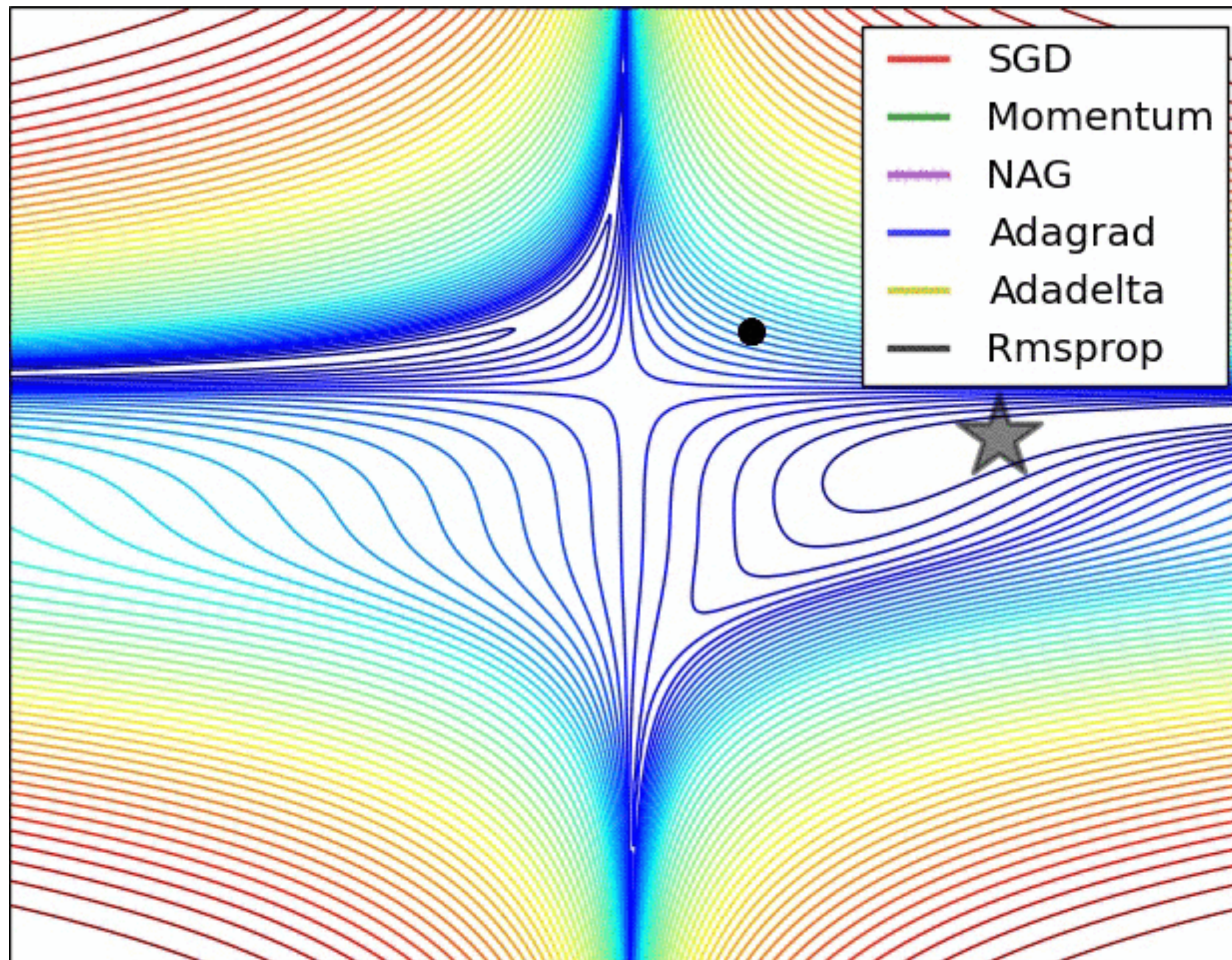
Common mini-batch sizes are 32/64/128 examples
e.g. Krizhevsky ILSVRC ConvNet used 256 examples

Fancier update formulas (momentum, Adagrad, RMSProp, Adam, …) — taught in 682

# The effects of different update form formulas



(image credits to Alec Radford)