

# Linear Classifiers

Last modified on 2022/05/07

Send typos and comments to [f.maire@qut.edu.au](mailto:f.maire@qut.edu.au)

Most of the following slides are adapted from  
*Fei-Fei Li and Andrej Karpathy*

# Outline

- Introductory example: image classification
- Nearest Neighbour classifier (review)
- L1 distance
- Linear classification
  - Losses (SVM loss, Cross-entropy loss)
  - Error landscapes
  - Training algorithms

# Image Classification: a core task in Computer Vision



X

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



cat

y

# The problem: semantic gap

Images are represented as  $R^d$   
arrays of numbers

- E.g.  $R^3$  with integers  
between [0, 255], where  
 $d=3$  represents 3 color  
channels (RGB)



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

What the computer sees

*The complex process of image processing: leading from raw images, through subsequent stages to the semantic interpretation of the image.*

# An image classifier

```
def predict(image):
    # *****
    return class_label
```

Unlike e.g. sorting a list of numbers,  
**no obvious way** to hard-code the algorithm for  
recognizing a cat, or other classes.

## Data-driven approach:

1. Collect a dataset of images and label them
2. Use Machine Learning to train an image classifier
3. Evaluate the classifier on a withheld set of test images

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

def predict(image):
    # evaluate the model on the image
    return class_label
```

Example training set



# First classifier: Nearest Neighbor Classifier

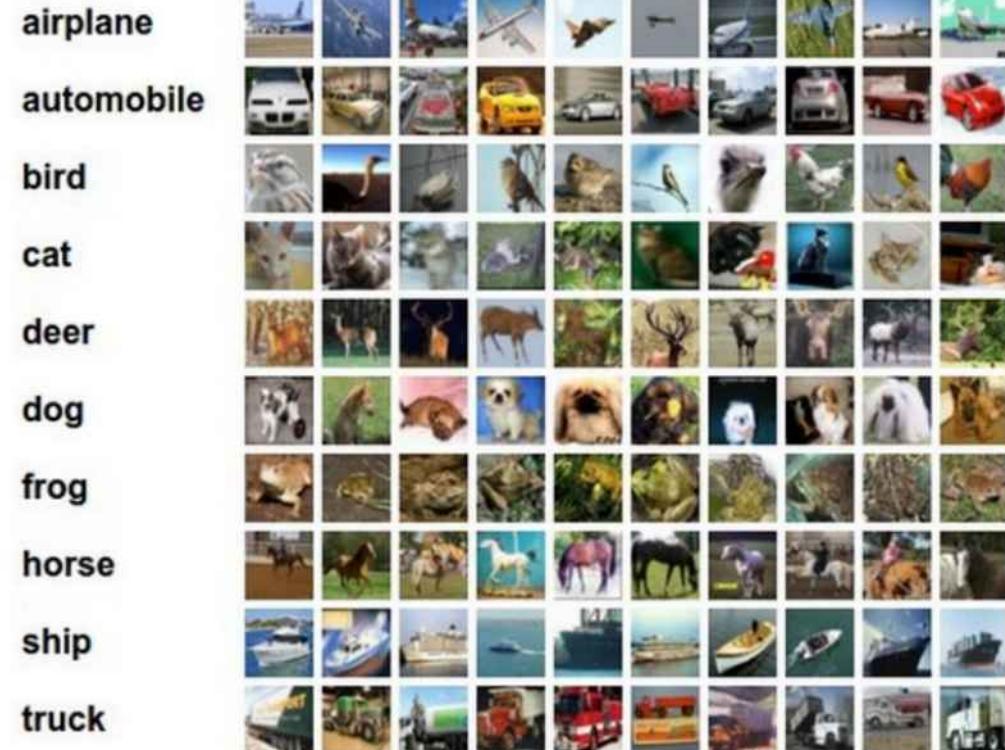
```
def train(train_images, train_labels):  
    # build a model of images -> labels  
  
def predict(image):  
    # evaluate the model on the image  
    return class_label
```



Remember all training images and their labels

Predict the label of the most similar training image

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



How do we compare the images? What is the **distance metric**?

**L1 distance:**

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

Where  $I_1$  denotes image 1,  
and  $p$  denotes each pixel

test image				training image				pixel-wise absolute value differences			
56	32	10	18	10	20	24	17	46	12	14	1
90	23	128	133	8	10	89	100	82	13	39	33
24	26	178	200	12	16	178	170	12	10	0	30
2	0	255	220	4	32	233	112	2	32	22	108

- = → 456

# Linear Classification

## 1. define a **score function**

(assume CIFAR-10 example so  
32 x 32 x 3 images, 10 classes)

class scores  
[10 x 1]

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

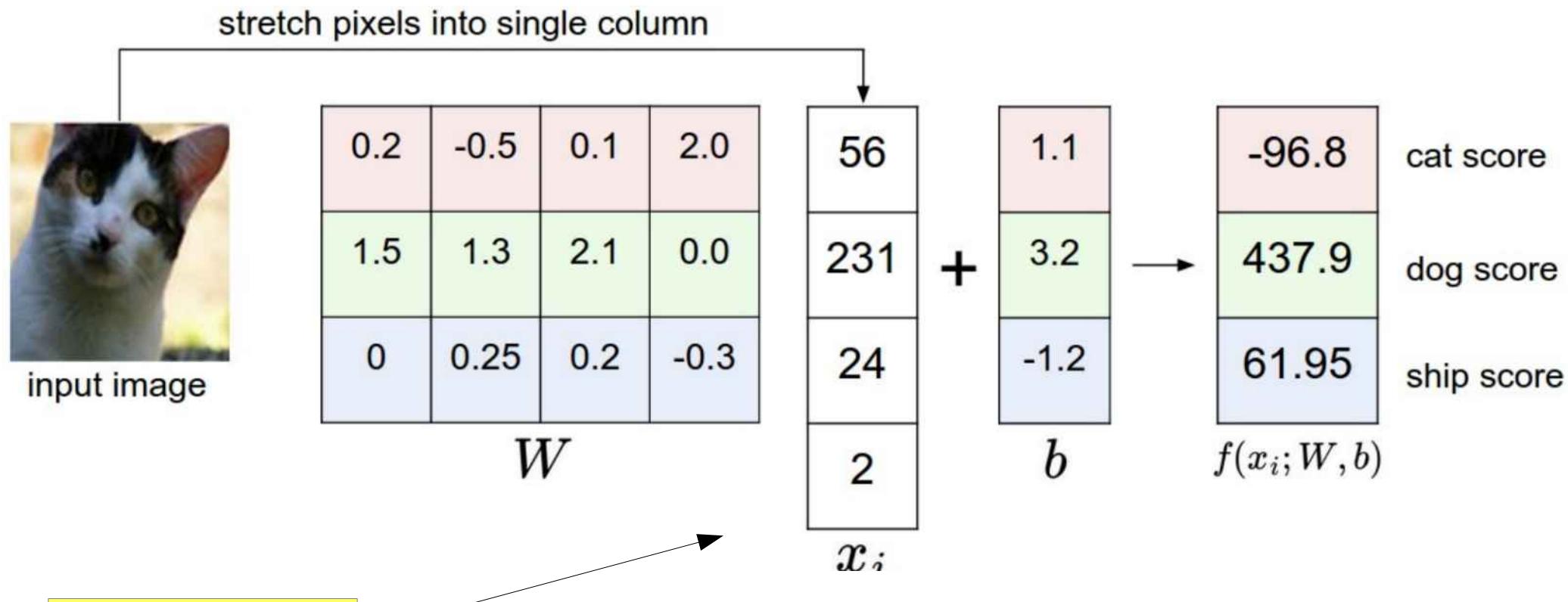
data (image)  
[3072 x 1]

weights  
[10 x 3072]

bias vector  
[10 x 1]

# Linear Classification

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



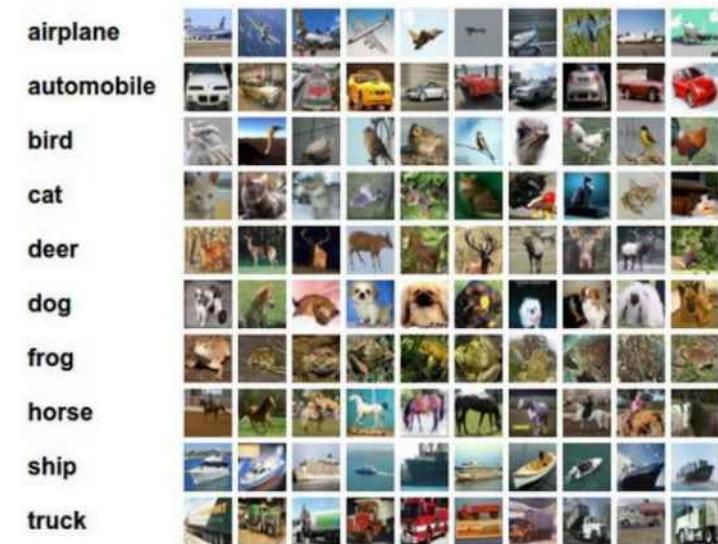
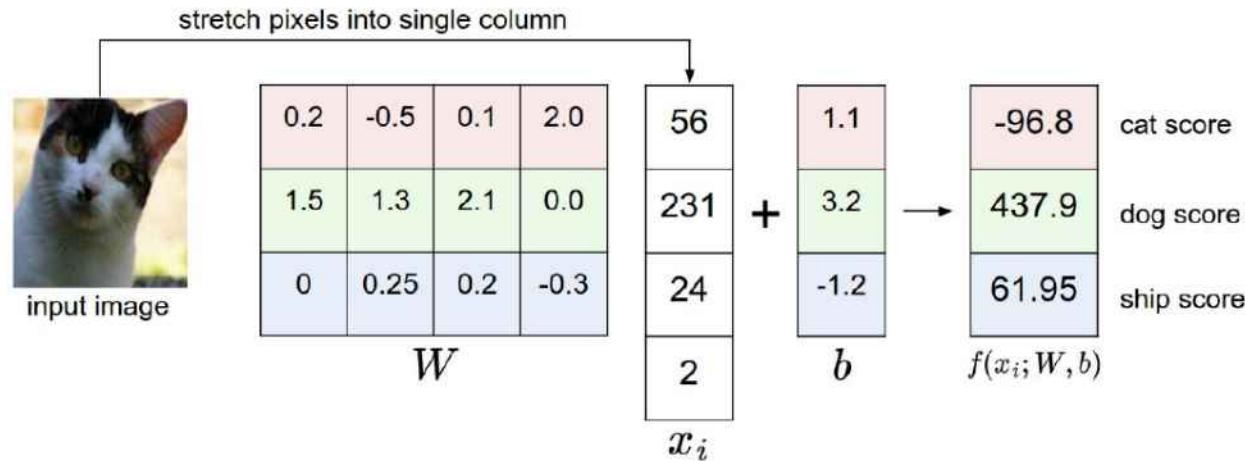
You have to imagine  
that the image has  
only 4 pixels!

# Interpreting a Linear Classifier

Question:

what can a linear classifier do?

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



# Interpreting a Linear Classifier

airplane  
automobile  
bird  
cat  
deer  
dog  
frog  
horse  
ship  
truck



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

Example training  
classifiers on CIFAR-10:

plane



car



bird



cat



deer



dog



frog



horse



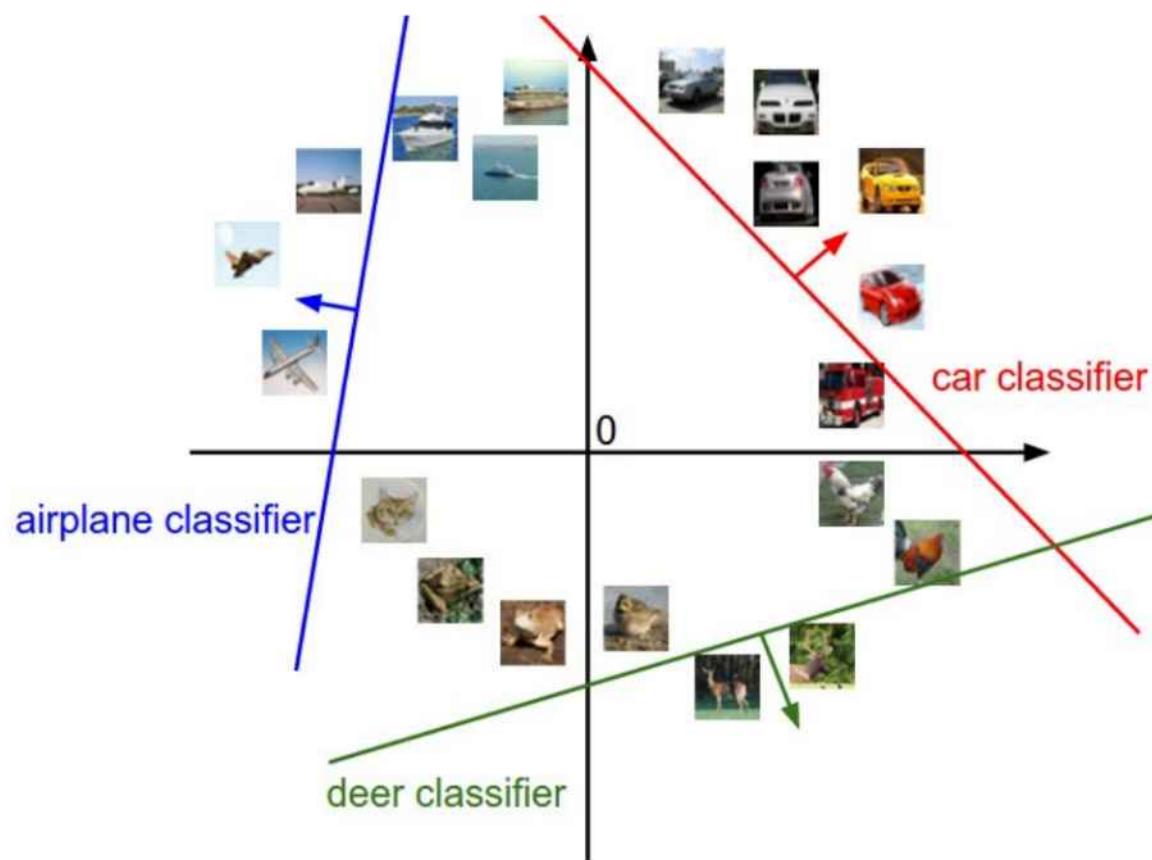
ship



truck



# Interpreting a Linear Classifier



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

# Bias trick

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

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

$W$

56	1.1
231	3.2
24	-1.2
2	

$x_i$

+



0.2	-0.5	0.1	2.0	1.1
1.5	1.3	2.1	0.0	3.2
0	0.25	0.2	-0.3	-1.2

$W$        $b$

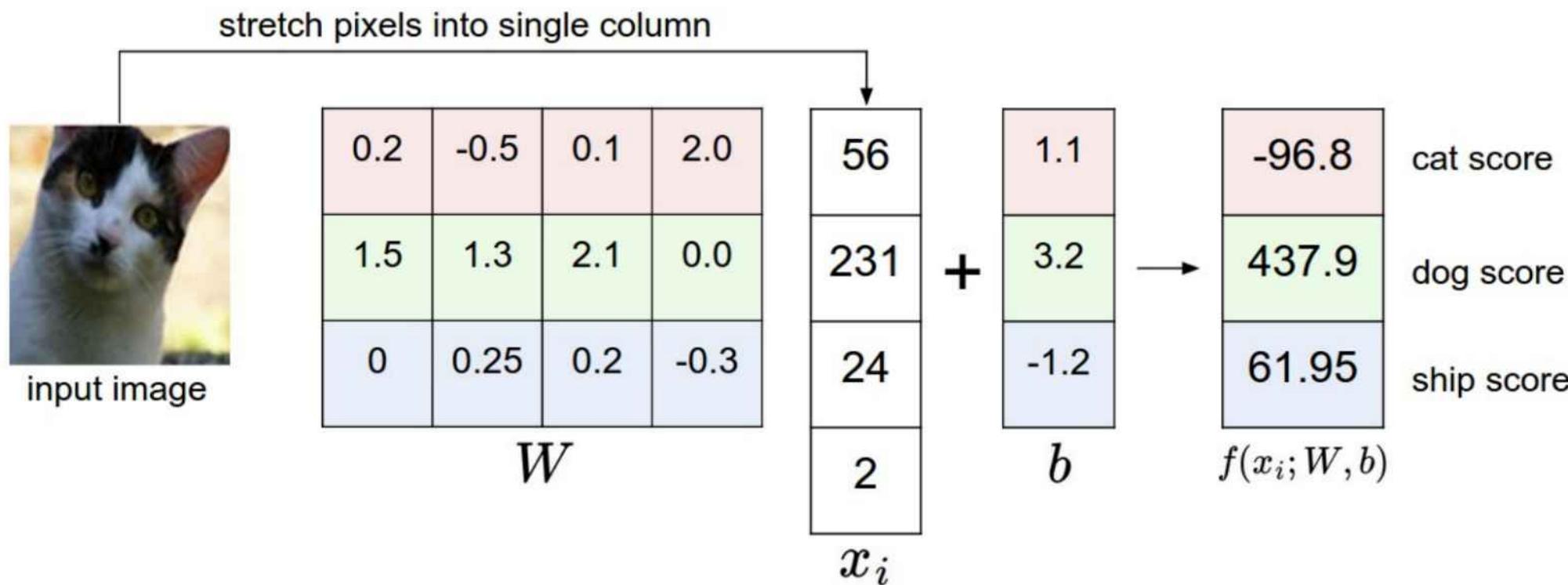
new, single  $W$

56
231
24
2
1

$x_i$

# So far:

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



## Define a loss function (or cost function, or objective)

- scores, label  $\longrightarrow$  loss.  
 $f(x_i, W)$        $y_i$                            $L_i$
- 

**Example:**

$$f(x_i, W) = [13, -7, 11]$$

$$y_i = 0$$

**Question:** if you were to assign a single number to how “unhappy” you are with these scores, what would you do?

Define a loss function (or cost function, or objective)  
One (of many ways) to do it: **Multiclass SVM Loss**

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

↑  
loss due to  
example i

↑  
sum over all  
incorrect labels

↑  
difference between the correct class  
score and incorrect class score



$$\text{Example: } L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

$$f(x_i, W) = [13, -7, 11]$$

$$y_i = 0$$

e.g. 10

---

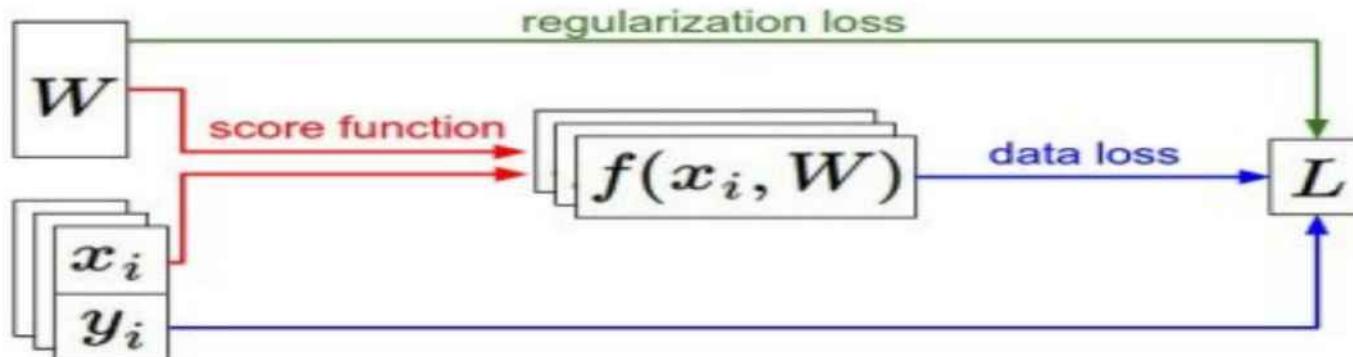
$$L_i = \max(0, -7 - 13 + 10) + \max(0, 11 - 13 + 10)$$

# L2 Regularization

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

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} \left[ \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta) \right] + \boxed{\lambda R(W)}$$

Regularization strength



## Softmax Classifier

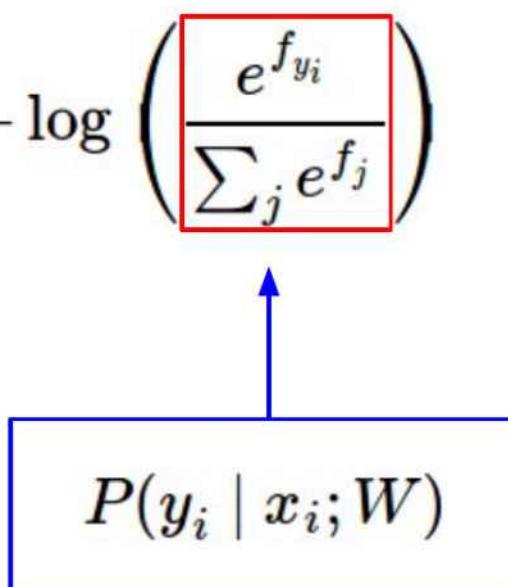
$$f(x_i, W) = Wx_i$$

score function  
is the same

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

softmax function

$P(y_i | x_i; W)$



i.e. we're minimizing  
the negative log  
likelihood.

$$[1, -2, 0] \rightarrow [e^1, e^{-2}, e^0] = [2.71, 0.14, 1] \rightarrow [0.7, 0.04, 0.26]$$

matrix multiply + bias offset

0.01	-0.05	0.1	0.05
0.7	0.2	0.05	0.16
0.0	-0.45	-0.2	0.03

$W$

-15	0.0
22	0.2
-44	-0.3
56	

+

$b$

$y_i$

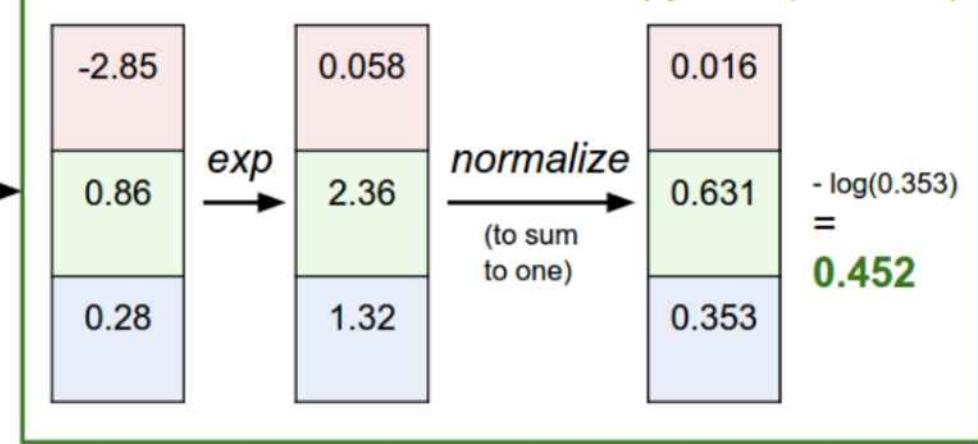
2

-2.85
0.86
0.28

hinge loss (SVM)

$$\max(0, -2.85 - 0.28 + 1) + \max(0, 0.86 - 0.28 + 1) = 1.58$$

cross-entropy loss (Softmax)



## Softmax vs. SVM

- Interpreting the probabilities from the Softmax

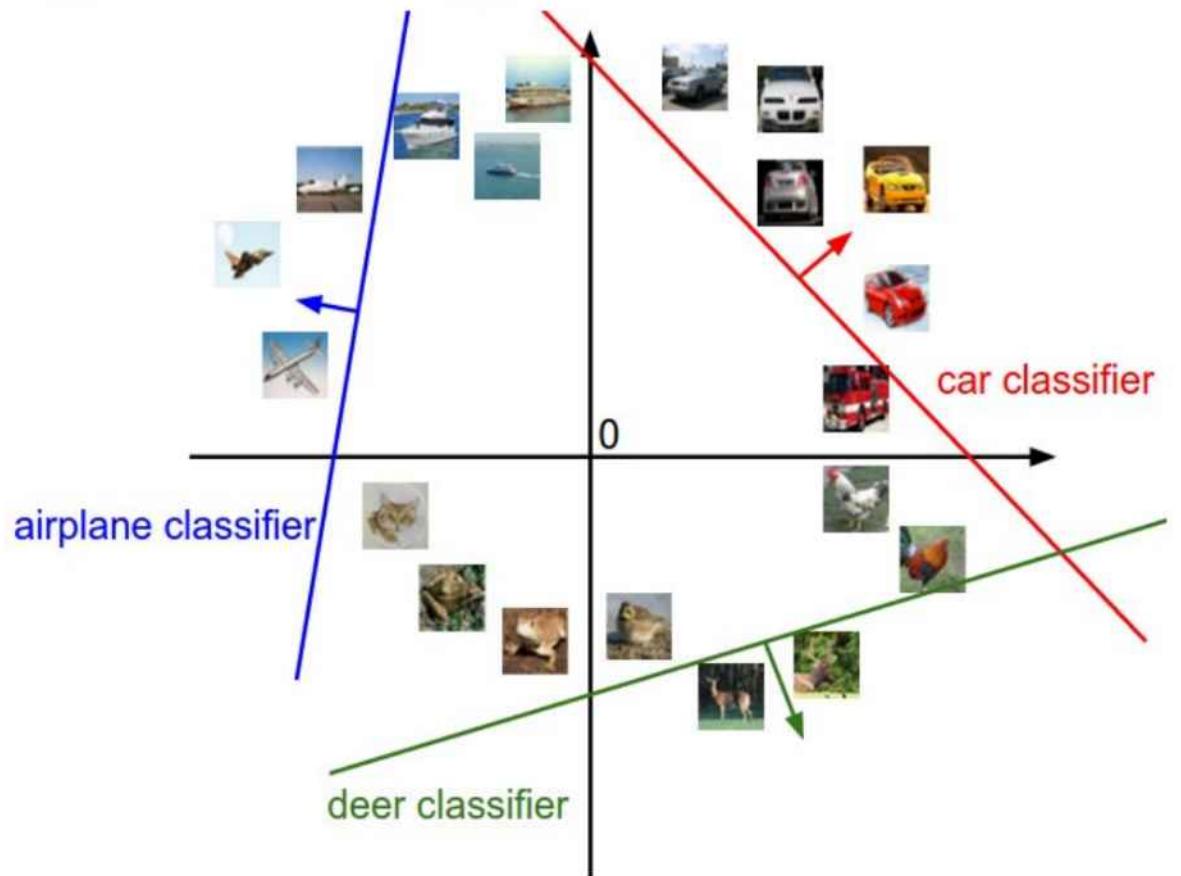
$$[1, -2, 0] \rightarrow [e^1, e^{-2}, e^0] = [2.71, 0.14, 1] \rightarrow [0.7, 0.04, 0.26]$$

suppose the weights W were only half as large:

$$[0.5, -1, 0] \rightarrow [e^{0.5}, e^{-1}, e^0] = [1.65, 0.37, 1] \rightarrow [0.55, 0.12, 0.33]$$

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + 1)$$

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



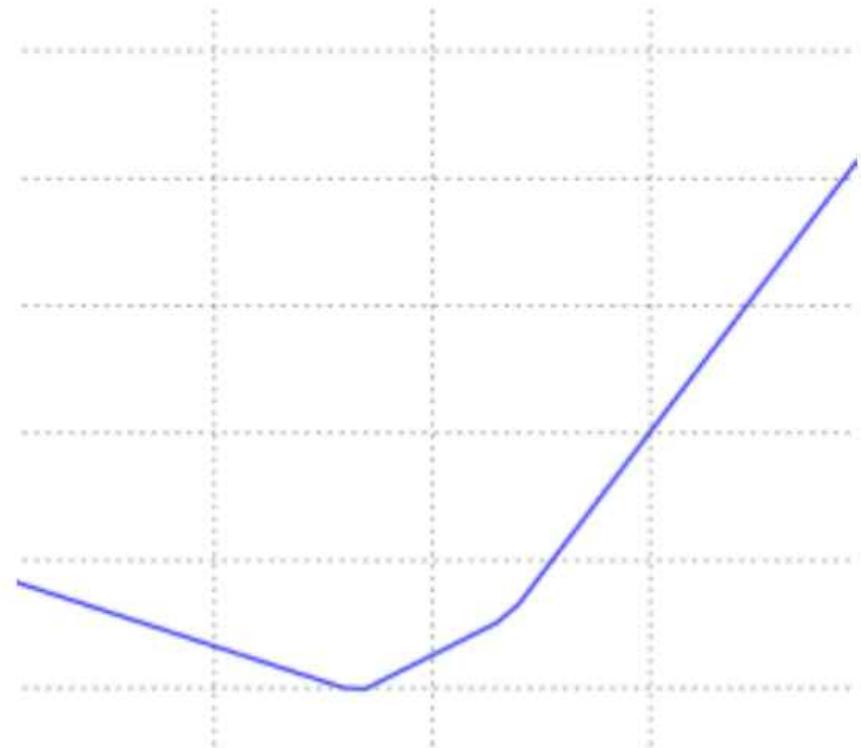
How do we find  $W$  and  $b$ ?

# Visualizing the (SVM) loss function

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

$$L_i(W + aW_1)$$

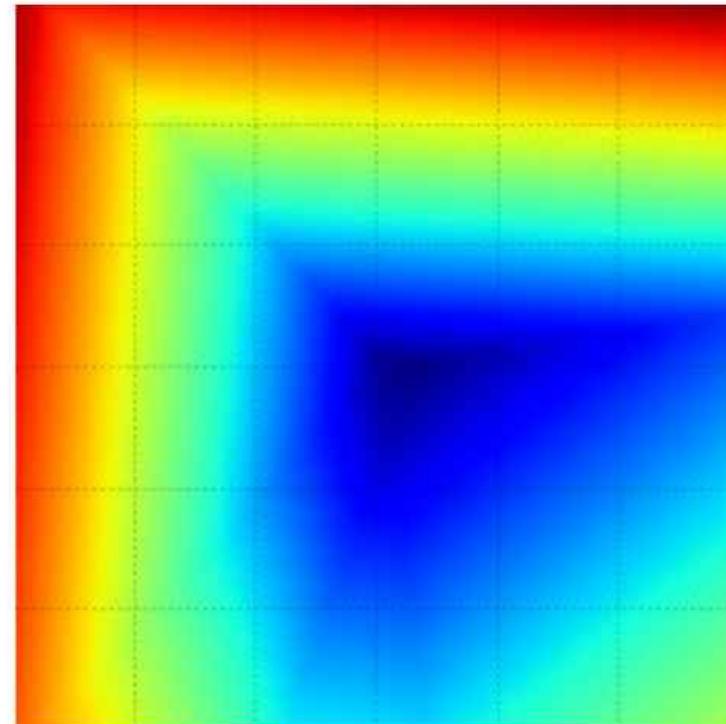
Plot of  $L_i$  versus  $a$



# Visualizing the (SVM) loss function

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

$$L_i(W + aW_1 + bW_2)$$



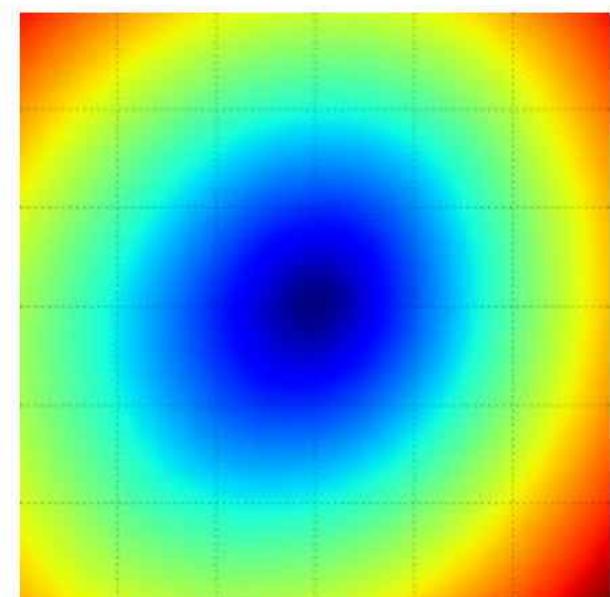
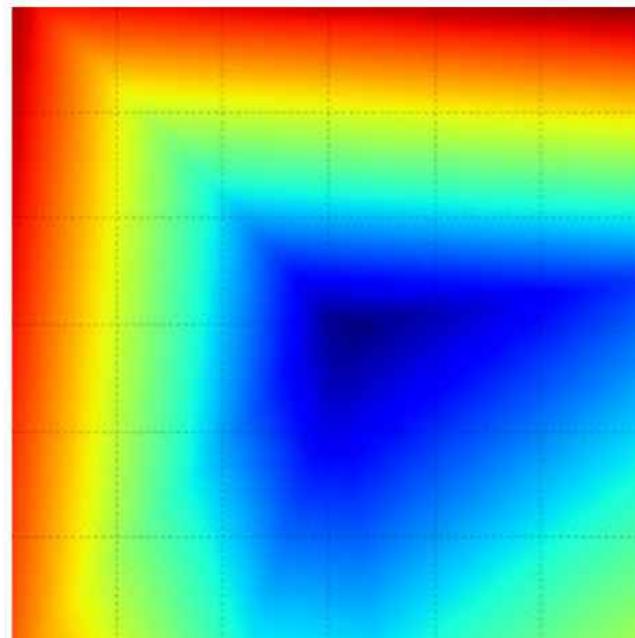
Plot of  $L_i$  versus  $a$  and  $b$

# Visualizing the (SVM) loss function

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

the full data loss:

$$L_i(W + aW_1 + bW_2)$$



# Visualizing the (SVM) loss function

$$L_i = \sum_{j \neq y_i} \left[ \max(0, w_j^T x_i - w_{y_i}^T x_i + 1) \right]$$

Suppose there are 3 examples with 3 classes (class 0, 1, 2 in sequence), then this becomes:

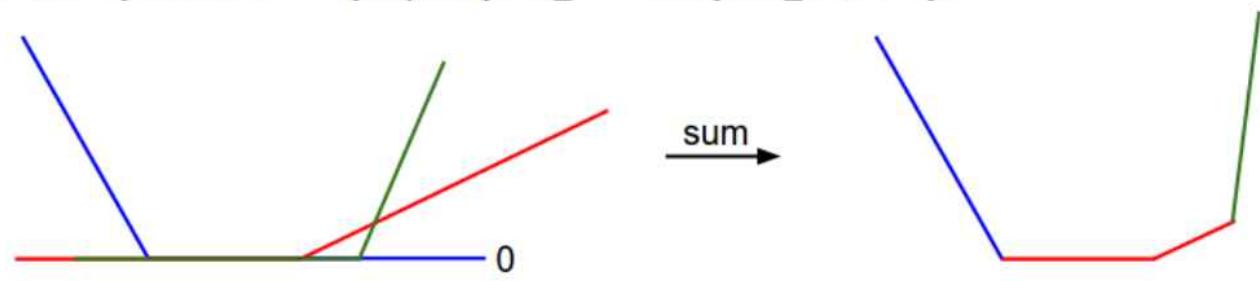
$$L_0 = \max(0, w_1^T x_0 - w_0^T x_0 + 1) + \max(0, w_2^T x_0 - w_0^T x_0 + 1)$$

$$L_1 = \max(0, w_0^T x_1 - w_1^T x_1 + 1) + \max(0, w_2^T x_1 - w_1^T x_1 + 1)$$

$$L_2 = \max(0, w_0^T x_2 - w_2^T x_2 + 1) + \max(0, w_1^T x_2 - w_2^T x_2 + 1)$$

$$L = (L_0 + L_1 + L_2)/3$$

One plot for each  $L_i$



sum →

Plot of  $L$

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

```
# 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
# ... (truncated: continues for 1000 lines)
```

## Lets see how well this works on the test set...

```
# 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
```

## Strategy #2: A better but still very bad idea solution: Random local search

```
W = np.random.randn(10, 3073) * 0.001 # generate random starting W
bestloss = float("inf")
for i in xrange(1000):
    step_size = 0.0001
    Wtry = W + np.random.randn(10, 3073) * step_size
    loss = L(Xtr_cols, Ytr, Wtry)
    if loss < bestloss:
        W = Wtry
        bestloss = loss
    print 'iter %d loss is %f' % (i, bestloss)
```

gives 21.4%!

## Strategy #3: Following the gradient

In 1-dimension, the derivative of a function:

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

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

# Evaluation the gradient numerically

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

in practice:

$$[f(x + h) - f(x - h)]/2h$$

“centered difference formula”

```
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) # evaluate 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
```

```
# to use the generic code above we want a function that takes a single argument
# (the weights in our case) so we close over X_train and Y_train
def CIFAR10_loss_fun(W):
    return L(X_train, Y_train, W)
```

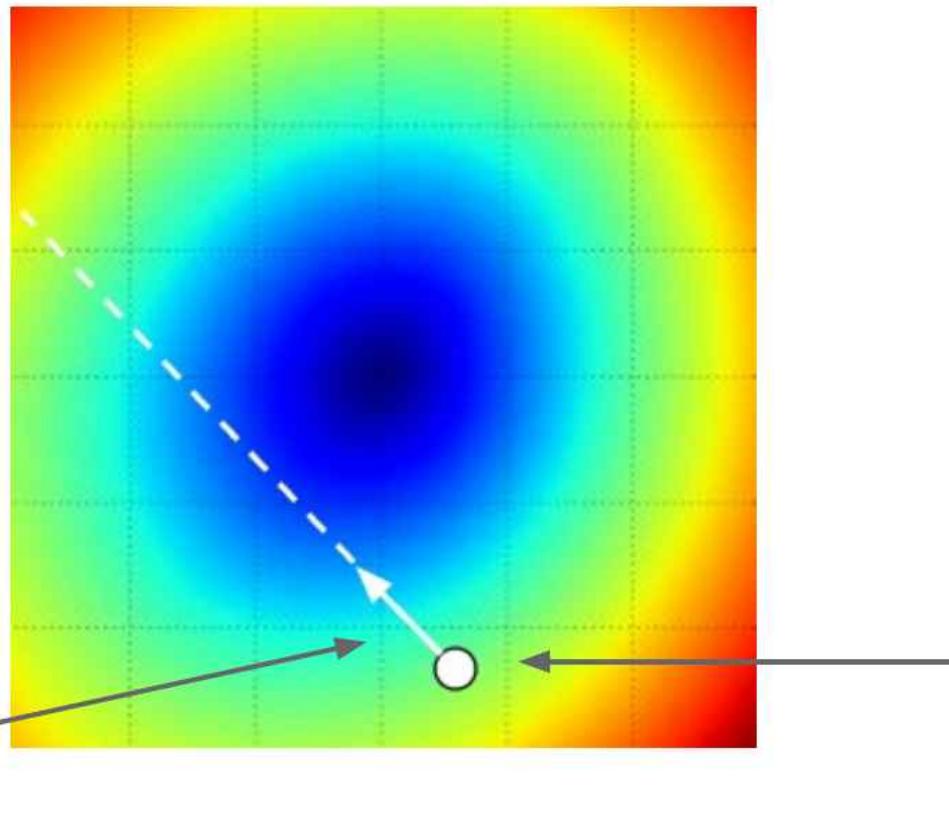
```
W = np.random.rand(10, 3073) * 0.001 # random weight vector
df = eval_numerical_gradient(CIFAR10_loss_fun, W) # get the gradient
```

# performing a parameter update

```
loss_original = CIFAR10_loss_fun(W) # the original loss
print 'original loss: %f' % (loss_original, )

# lets see the effect of multiple step sizes
for step_size_log in [-5,-4,-3,-2,-1,0,1,2]:
    step_size = 10 ** step_size_log
    W_new = W - step_size * df # new position in the weight space
    loss_new = CIFAR10_loss_fun(W_new)
    print 'for step size %f new loss: %f' % (step_size, loss_new)

# prints:
# original loss: 2.200718
# for step size 1.000000e-10 new loss: 2.200652
# for step size 1.000000e-09 new loss: 2.200057
# for step size 1.000000e-08 new loss: 2.194116
# for step size 1.000000e-07 new loss: 2.135493
# for step size 1.000000e-06 new loss: 1.647802
# for step size 1.000000e-05 new loss: 2.844355
# for step size 1.000000e-04 new loss: 25.558142
# for step size 1.000000e-03 new loss: 254.086573
# for step size 1.000000e-02 new loss: 2539.370888
# for step size 1.000000e-01 new loss: 25392.214036
```



original W

negative gradient direction

# The problems with numerical gradient:

- approximate
- very slow to evaluate

```
# to use the generic code above we want a function that takes a single argument
# (the weights in our case) so we close over X_train and Y_train
def CIFAR10_loss_fun(W):
    return L(X_train, Y_train, W)

W = np.random.rand(10, 3073) * 0.001 # random weight vector
df = eval_numerical_gradient(CIFAR10_loss_fun, W) # get the gradient
```

```
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
        x[ix] += h # increment by h
        fxh = f(x) # evaluate f(x + h)
        x[ix] -= h # 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
```

We need something better...

$$L_i = \sum_{j \neq y_i} \left[ \max(0, w_j^T x_i - w_{y_i}^T x_i + 1) \right]$$

$$\begin{aligned} L_i &= \max (0, w_1^T x_i - w_0^T x_i + 1) \\ &+ \max (0, w_2^T x_i - w_0^T x_i + 1) \\ &+ \max (0, w_3^T x_i - w_0^T x_i + 1) \end{aligned}$$



$$\nabla_{w_j} L_i = 1(w_j^T x_i - w_{y_i}^T x_i + \Delta > 0) x_i$$

$$\nabla_{w_{y_i}} L_i = - \left( \sum_{j \neq y_i} 1(w_j^T x_i - w_{y_i}^T x_i + \Delta > 0) \right) x_i$$

- 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
```

# 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 ~100 examples.  
e.g. Krizhevsky ILSVRC ConvNet used 256 examples

# Stochastic Gradient Descent (SGD)

- use a single example at a time

```
# 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
```

1



(also sometimes called **on-line** Gradient Descent)

- Always use mini-batch gradient descent
- Incorrectly refer to it as “doing SGD” as everyone else  
(or call it batch gradient descent)
- The mini-batch size is a hyperparameter, but it is not very common to cross-validate over it (usually based on practical concerns, e.g. space/time efficiency)

# The dynamics of Gradient Descent

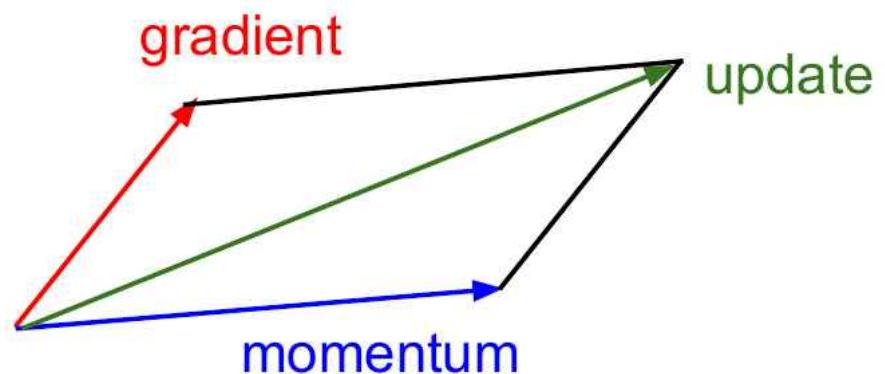
pull some weights up and some down

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} \left[ \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta) \right] + \lambda \sum_k \sum_l W_{k,l}^2$$

$$L = \frac{1}{N} \sum_i -\log \left( \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right) + \lambda \sum_k \sum_l W_{k,l}^2$$

always pull the weights down

# Momentum Update



```
weights_grad = evaluate_gradient(loss_fun, data, weights)
vel = vel * 0.9 - step_size * weights_grad
weights += vel
```