



# Машинное зрение

---

Лекция 3. Классификация  
изображений

1. Линейные классификаторы.
2. SVM классификатор. Hinge loss.
3. Softmax классификатор. Cross entropy loss.
4. Оптимизация параметров.

Материалы для более глубокого изучения:

- [cs231n.stanford.edu](http://cs231n.stanford.edu) (Convolution NN)
- <http://www.machinelearning.ru/wiki/images/archive/a/a0/2015031617222%21Voron-ML-Lin-SVM.pdf> (Лекция Воронцова К.В. ШАД)

Сроки сдачи CodeLab\_1:

- 26 ноября 2021 г.

# Image Classification: A core task in Computer Vision

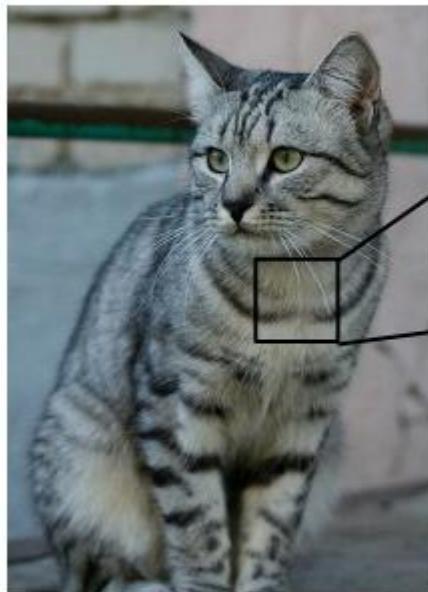


This image by [Nikita](#) is  
licensed under [CC-BY 2.0](#)



cat

## The Problem: Semantic Gap



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

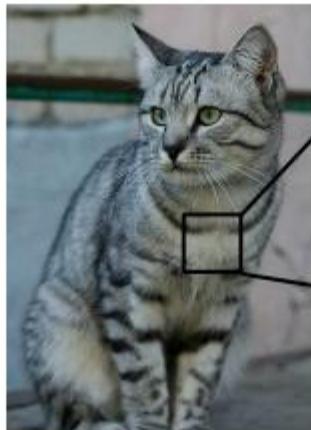
[[105 112 108 111 104 99 106 90 95 103 112 119 184 97 93 87]  
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]  
[ 76 85 90 105 128 185 87 96 95 99 115 112 106 103 99 85]  
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]  
[ 106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]  
[ 114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 96 95]  
[ 128 137 147 183 65 83 88 80 65 52 54 74 84 102 93 85 82]  
[ 133 137 147 183 65 83 88 80 65 52 54 74 84 102 93 85 82]  
[ 128 137 144 140 109 95 86 78 62 65 63 63 68 73 86 101]  
[ 125 133 148 137 119 121 117 94 65 79 80 65 54 64 72 98]  
[ 127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]  
[ 115 114 109 123 158 148 131 118 113 109 109 92 74 65 72 78]  
[ 89 93 98 97 108 147 131 118 113 114 113 109 106 95 77 80]  
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]  
[ 62 65 82 89 78 71 88 101 124 126 118 101 107 114 131 119]  
[ 63 65 75 88 89 73 62 81 128 138 135 185 81 98 118 118]  
[ 87 65 71 87 106 95 60 45 76 138 128 187 92 94 105 112]  
[ 118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]  
[ 164 146 112 80 82 120 124 104 76 48 45 66 88 101 102 109]  
[ 157 170 157 120 93 86 114 132 112 97 69 55 70 82 99 94]  
[ 130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]  
[ 128 112 96 117 150 144 120 115 104 107 102 93 87 81 72 79]  
[ 123 107 96 86 93 112 153 149 122 109 104 75 88 107 112 99]  
[ 122 121 102 80 82 86 94 117 145 148 153 182 58 78 92 107]  
[ 122 164 148 103 71 56 78 83 93 103 119 139 102 61 60 84]]

What the computer sees

An image is a tensor of integers  
between [0, 255]:

e.g. 800 x 600 x 3  
(3 channels RGB)

## Challenges: Viewpoint variation



1189	112	148	111	184	98	146	93	98	185	125	139	184	97	93	971		
1	93	98	182	188	184	79	86	182	89	185	123	136	129	185	94	851	
1	78	85	98	188	128	185	87	98	85	85	135	132	185	185	93	851	
1	68	85	85	93	128	131	127	188	85	89	182	89	96	93	181	941	
1	186	95	63	64	88	92	88	85	185	187	189	88	75	84	96	951	
1	141	138	137	142	143	144	145	146	147	148	149	150	151	152	153	1511	
1	133	137	142	189	85	85	88	85	82	83	78	80	182	93	85	821	
1	128	137	144	148	185	85	86	79	82	65	83	83	69	73	95	1811	
1	128	139	148	137	131	125	137	84	88	79	88	85	54	84	72	881	
1	123	131	138	139	141	142	143	144	145	146	147	148	149	150	151	1511	
1	128	134	189	129	139	148	132	138	133	189	189	89	74	88	72	781	
1	89	93	98	97	189	147	135	138	133	154	133	189	185	95	77	801	
1	63	77	88	85	85	77	79	182	128	137	129	127	128	128	128	871	
1	62	65	62	69	78	75	68	181	124	128	129	185	187	114	131	1191	
1	64	65	65	65	65	65	65	65	65	65	65	65	65	65	65	651	
1	87	85	73	87	186	95	69	45	76	139	125	187	92	84	195	3121	
1	118	89	92	86	117	123	126	66	45	31	85	84	89	85	182	3871	
1	164	146	112	88	82	128	124	184	78	49	45	86	89	181	182	3891	
1	153	123	122	122	122	122	122	122	122	122	122	122	122	122	122	1221	
1	128	128	124	126	129	148	189	128	125	124	124	87	85	23	89	801	
1	128	128	124	126	129	148	189	128	125	124	124	87	85	23	89	801	
1	128	122	96	117	158	144	128	115	184	187	182	83	87	81	72	791	
1	129	187	98	88	85	93	112	153	149	132	189	184	75	89	187	152	991
1	122	123	182	88	82	85	94	117	145	149	153	182	58	78	92	1871	
1	122	164	148	149	71	58	78	83	92	183	129	129	182	81	85	8411	

All pixels change when  
the camera moves!



## Challenges: Background Clutter

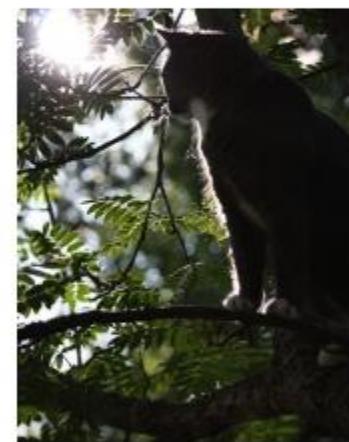


[This image](#) is CC0 1.0 public domain



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

## Challenges: Illumination



## Challenges: Occlusion



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



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



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

## Challenges: Deformation



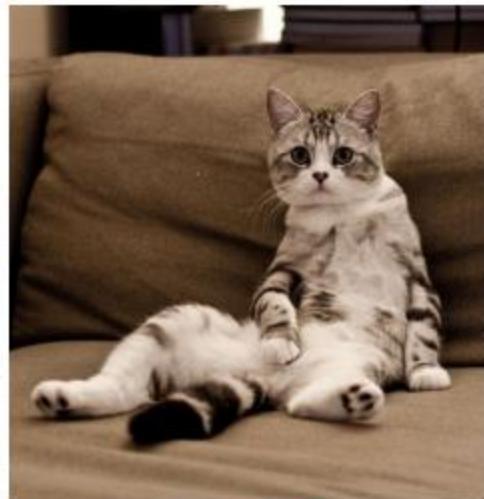
[This image](#) by [Umberto Salvagnin](#)  
is licensed under [CC-BY 2.0](#)



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



[This image](#) by [sare\\_bear](#) is  
licensed under [CC-BY 2.0](#)



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

## Challenges: Intraclass variation



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

### An image classifier

```
def classify_image(image):  
    # Some magic here?  
    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.

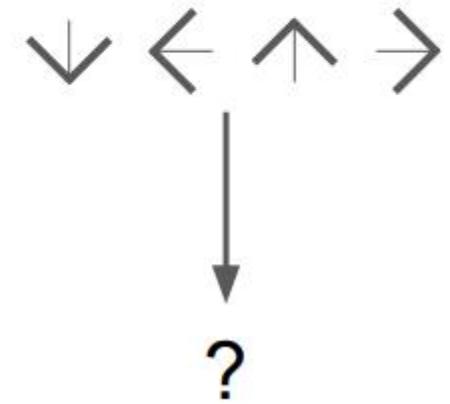
## Attempts have been made



Find edges



Find corners

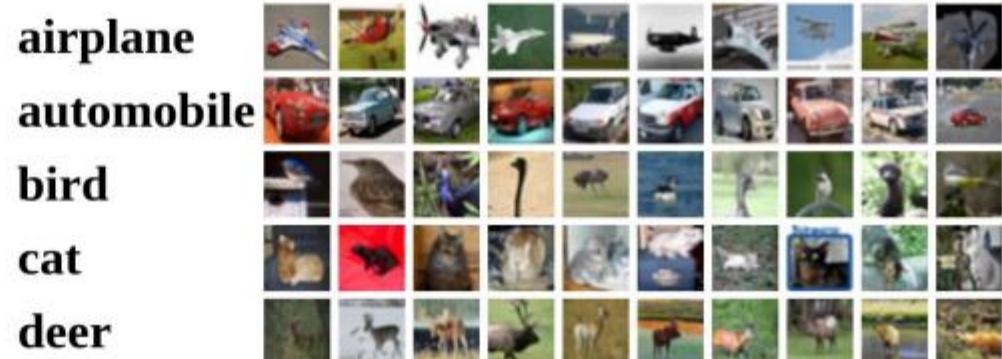


## Machine Learning: Data-Driven Approach

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

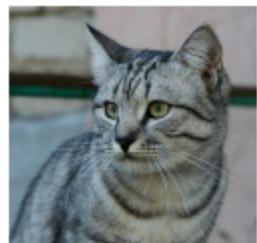
Example training set

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



## Parametric Approach

Image



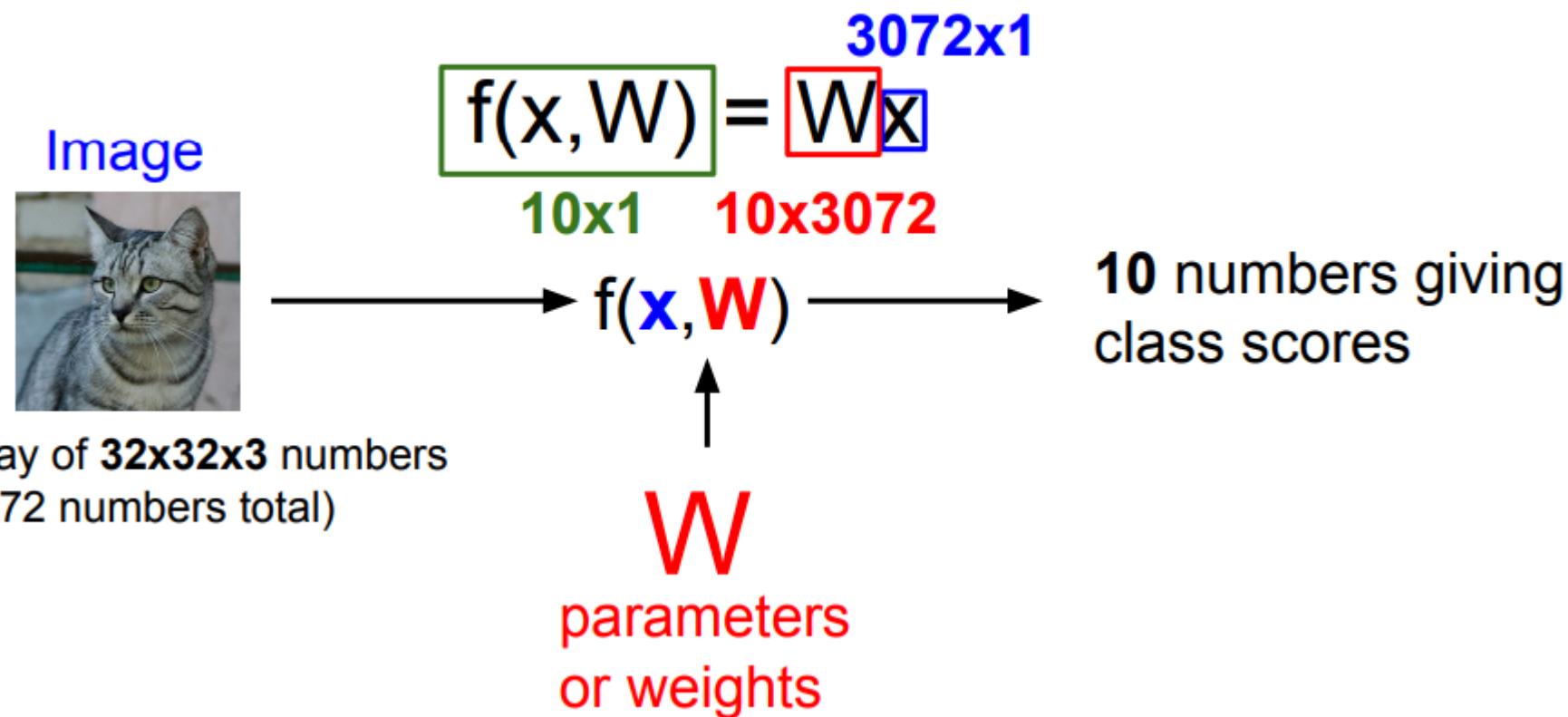
Array of **32x32x3** numbers  
(3072 numbers total)

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

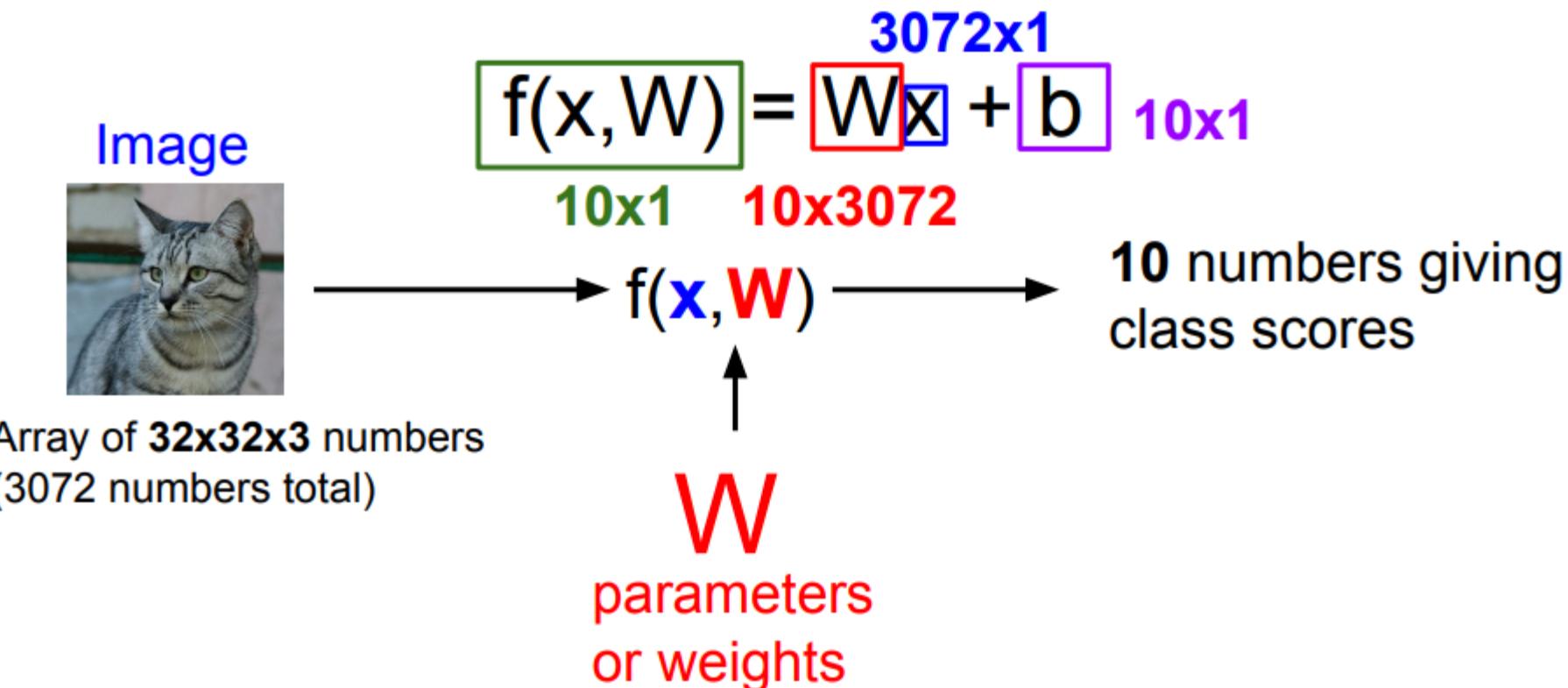
**W**  
parameters  
or weights

**10** numbers giving  
class scores

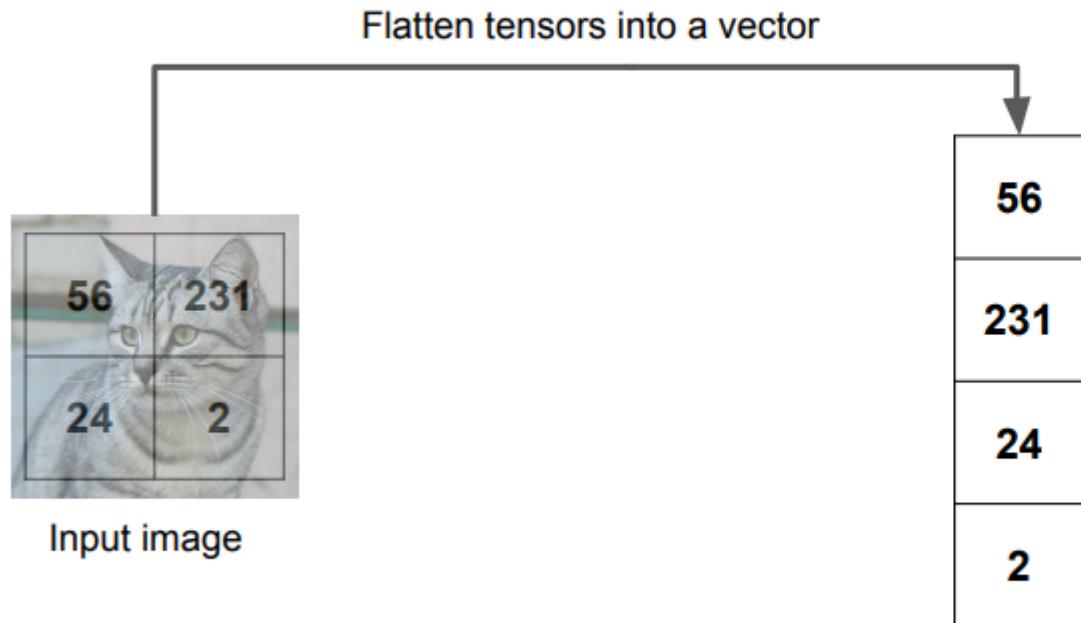
## Parametric Approach: Linear Classifier



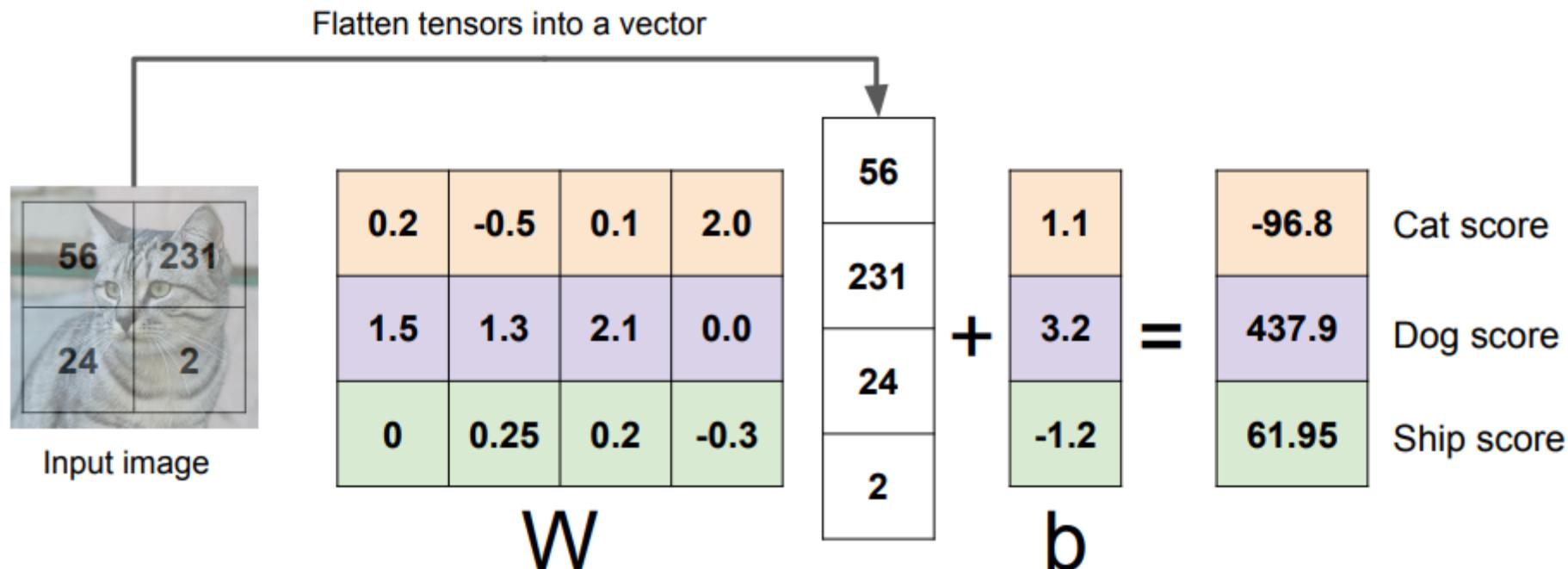
## Parametric Approach: Linear Classifier



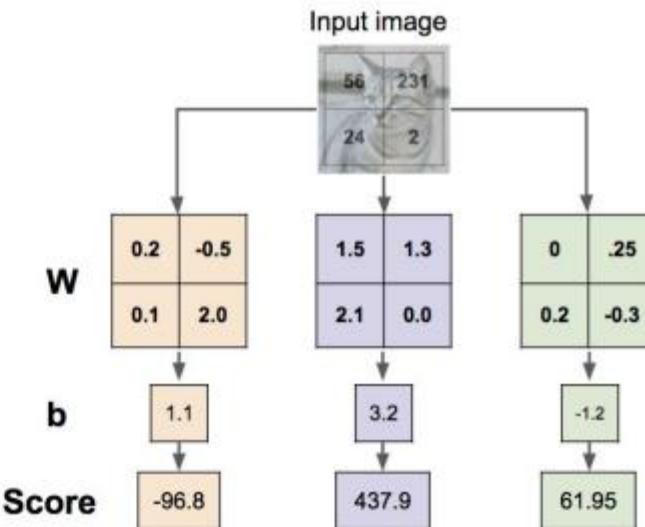
Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)



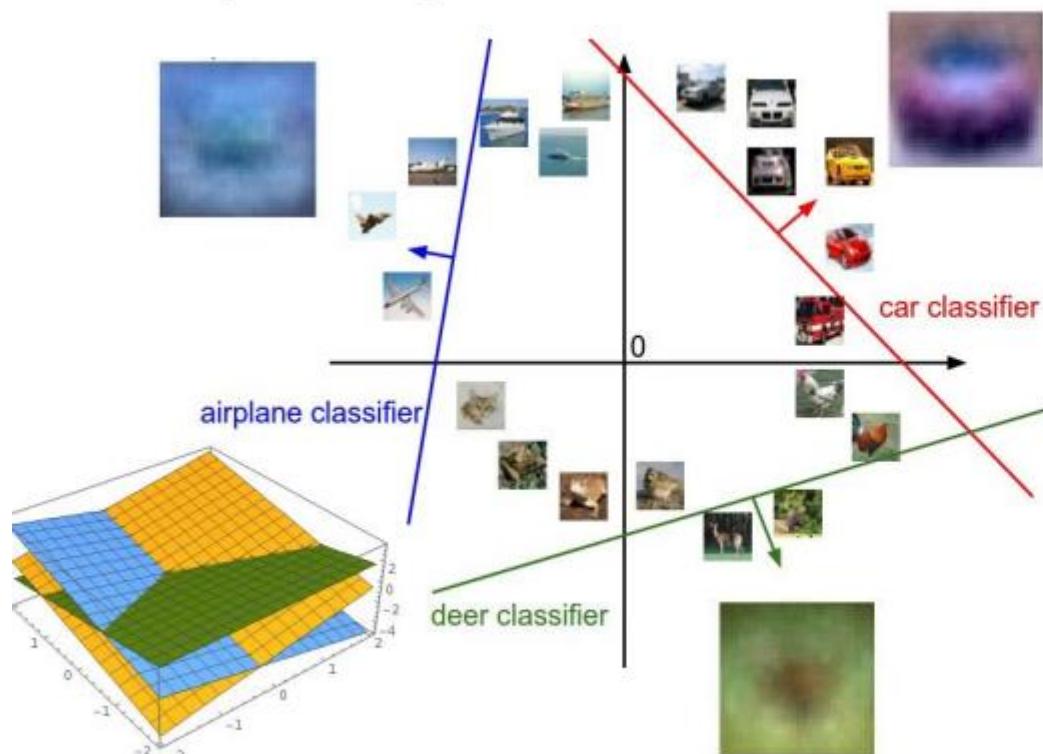
Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)



## Interpreting a Linear Classifier: Visual Viewpoint



## Interpreting a Linear Classifier: Geometric Viewpoint



$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + b$$



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

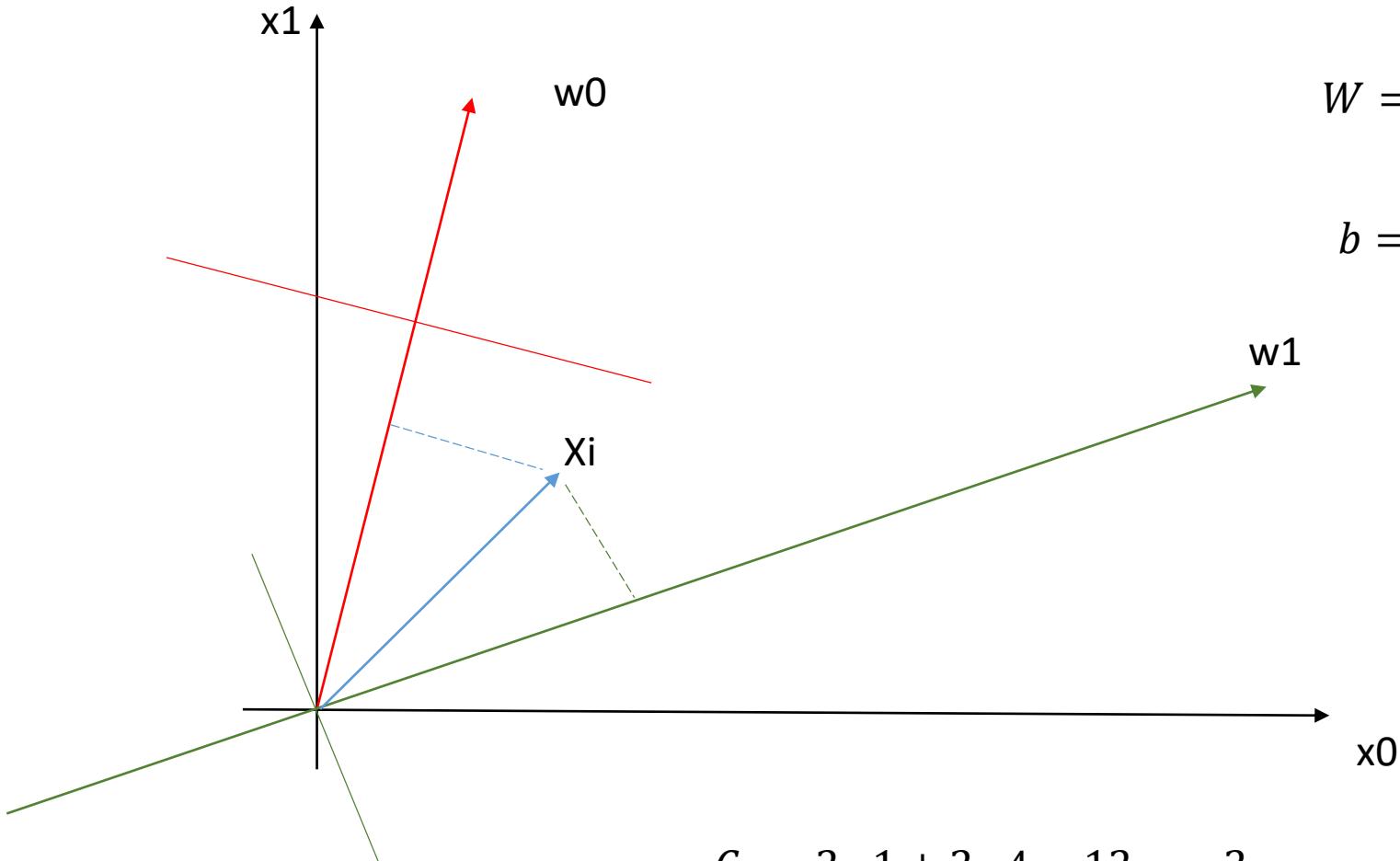
# Линейные классификаторы

<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>

$$X_i = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$W = \begin{bmatrix} 1 & 4 \\ 5 & 3 \end{bmatrix}$$

$$b = \begin{bmatrix} -12 \\ 0 \end{bmatrix}$$



$$C_0 = 2 \cdot 1 + 2 \cdot 4 - 12 = -2$$
$$C_1 = 2 \cdot 5 + 2 \cdot 3 + 0 = 16$$

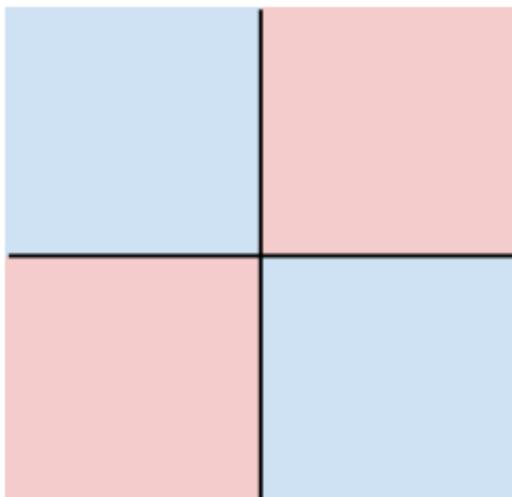
## Hard cases for a linear classifier

**Class 1:**

First and third quadrants

**Class 2:**

Second and fourth quadrants

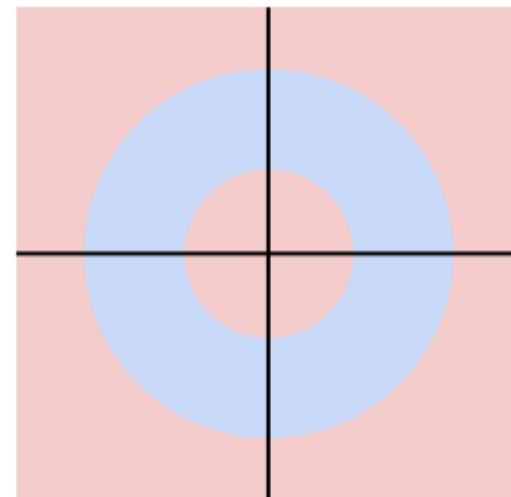


**Class 1:**

$1 \leq L_2 \text{ norm} \leq 2$

**Class 2:**

Everything else

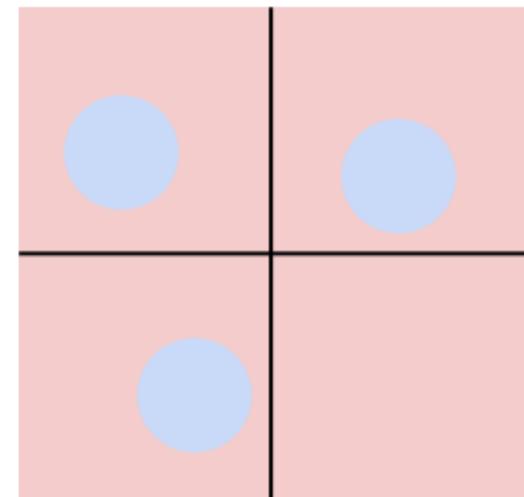


**Class 1:**

Three modes

**Class 2:**

Everything else



## Support vector machine

### 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)**



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

## SVM классификатор. Hinge loss.

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



cat	<b>3.2</b>	1.3	2.2
car	<b>5.1</b>	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

## SVM классификатор. Hinge loss.

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

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

Given a dataset of examples

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

Where  $x_i$  is image and  
 $y_i$  is (integer) label

## SVM классификатор. Hinge loss.

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



cat	<b>3.2</b>	1.3	2.2
car	<b>5.1</b>	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

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

Given a dataset of examples

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

Where  $x_i$  is image and  
 $y_i$  is (integer) label

Loss over the dataset is a average of loss over examples:

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

## SVM классификатор. Hinge loss.

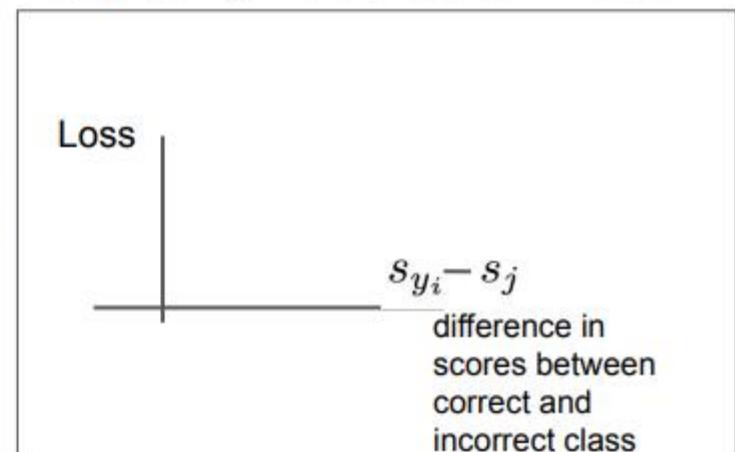
Suppose: 3 training examples, 3 classes.

With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

### Interpreting Multiclass SVM loss:



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

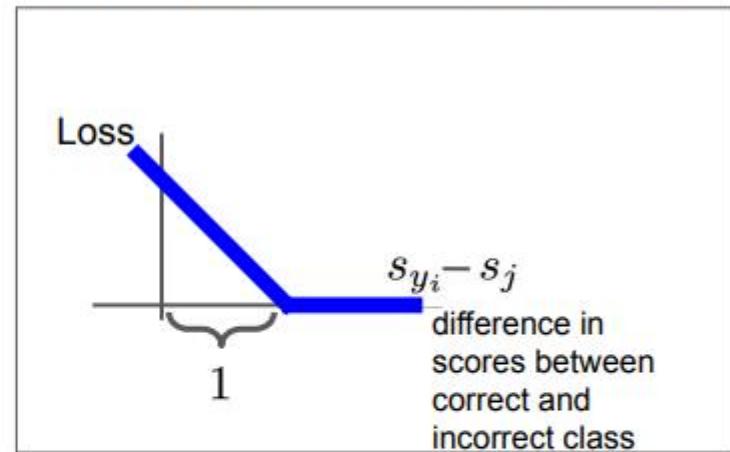
## SVM классификатор. Hinge loss.

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

### Interpreting Multiclass SVM loss:



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

## SVM классификатор. Hinge loss.

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		

### Multiclass SVM loss:

Given an example  $(x_i, y_i)$   
where  $x_i$  is the image and  
where  $y_i$  is the (integer) label,

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

the SVM loss has the form:

$$\begin{aligned}L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\&= \max(0, 5.1 - 3.2 + 1) \\&\quad + \max(0, -1.7 - 3.2 + 1) \\&= \max(0, 2.9) + \max(0, -3.9) \\&= 2.9 + 0 \\&= 2.9\end{aligned}$$

## SVM классификатор. Hinge loss.

Suppose: 3 training examples, 3 classes.

With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
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:  $s = f(x_i, W)$

the SVM loss has the form:

$$\begin{aligned}L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\&= \max(0, 1.3 - 4.9 + 1) \\&\quad + \max(0, 2.0 - 4.9 + 1) \\&= \max(0, -2.6) + \max(0, -1.9) \\&= 0 + 0 \\&= 0\end{aligned}$$

## SVM классификатор. Hinge loss.

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



cat	<b>3.2</b>	1.3	<b>2.2</b>
car	<b>5.1</b>	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	2.9	0	<b>12.9</b>

### Multiclass SVM loss:

Given an example  $(x_i, y_i)$   
where  $x_i$  is the image and  
where  $y_i$  is the (integer) label,

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

the SVM loss has the form:

$$\begin{aligned}L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\&= \max(0, 2.2 - (-3.1) + 1) \\&\quad + \max(0, 2.5 - (-3.1) + 1) \\&= \max(0, 6.3) + \max(0, 6.6) \\&= 6.3 + 6.6 \\&= 12.9\end{aligned}$$

Suppose: 3 training examples, 3 classes.

With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	<b>2.9</b>	<b>0</b>	<b>12.9</b>

## Multiclass SVM loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and where  $y_i$  is the (integer) label,

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

the SVM loss has the form:

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

Loss over full dataset is average:

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

$$\begin{aligned} L &= (2.9 + 0 + 12.9)/3 \\ &= 5.27 \end{aligned}$$

Suppose: 3 training examples, 3 classes.

With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	1.3
car	4.9
frog	2.0
Losses:	0

## Multiclass SVM loss:

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

Q1: What happens to loss if car scores decrease by 0.5 for this training example?

Q2: what is the min/max possible SVM loss  $L_i$ ?

Q3: At initialization  $W$  is small so all  $s \approx 0$ . What is the loss  $L_i$ , assuming  $N$  examples and  $C$  classes?

Suppose: 3 training examples, 3 classes.

With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
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:  $s = f(x_i, W)$

the SVM loss has the form:

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

Q5: What if we used mean instead of sum?

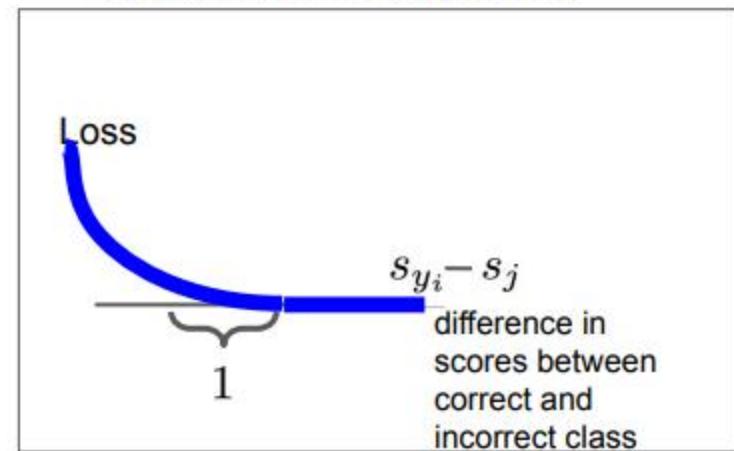
## SVM классификатор. Hinge loss.

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	<b>2.9</b>	<b>0</b>	<b>12.9</b>

### Multiclass SVM loss:



Q6: What if we used

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

## Multiclass SVM Loss: Example code

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

```
def L_i_vectorized(x, y, W):
    scores = W.dot(x)
    margins = np.maximum(0, scores - scores[y] + 1) # First calculate scores
    margins[y] = 0 # Then calculate the margins  $s_j - s_{y_i} + 1$ 
    loss_i = np.sum(margins) # only sum j is not  $y_i$ , so when  $j = y_i$ , set to zero.
    return loss_i # sum across all j
```

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

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

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

**No!  $2W$  is also has  $L = 0!$**

## SVM классификатор. Hinge loss

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	2.9	<b>0</b>	

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

**Before:**

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

**With  $W$  twice as large:**

$$\begin{aligned} &= \max(0, 2.6 - 9.8 + 1) \\ &\quad + \max(0, 4.0 - 9.8 + 1) \\ &= \max(0, -6.2) + \max(0, -4.8) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

# Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions}} + \lambda R(W)$$


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

**Regularization:** Prevent the model from doing *too well* on training data

# Regularization

$\lambda$  = regularization strength  
(hyperparameter)

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \lambda R(W)$$

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

**Regularization:** Prevent the model from doing *too well* on training data

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

## Regularization: Expressing Preferences

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

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

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

L2 Regularization

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

Which of w1 or w2 will  
the L2 regularizer prefer?

L2 regularization likes to  
“spread out” the weights

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

### Softmax Classifier (Multinomial Logistic Regression)

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



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

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

Softmax  
Function

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

## Softmax Classifier (Multinomial Logistic Regression)



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

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

Probabilities  
must be  $\geq 0$

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

Softmax  
Function

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

unnormalized probabilities

probabilities

## Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Softmax Function

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

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

cat  
car  
frog

<b>3.2</b>
5.1
-1.7

Unnormalized  
log-probabilities / logits

$\exp$

<b>24.5</b>
164.0
0.18

unnormalized  
probabilities

normalize

<b>0.13</b>
0.87
0.00

probabilities

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

**Maximum Likelihood Estimation**  
Choose weights to maximize the likelihood of the observed data  
(See CS 229 for details)

## Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Softmax Function

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

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

cat  
car  
frog

$$\begin{matrix} 3.2 \\ 5.1 \\ -1.7 \end{matrix}$$

Unnormalized  
log-probabilities / logits

$\exp$

$$\begin{matrix} 24.5 \\ 164.0 \\ 0.18 \end{matrix}$$

unnormalized  
probabilities

normalize

$$\begin{matrix} 0.13 \\ 0.87 \\ 0.00 \end{matrix}$$

probabilities

→ compare ←

Kullback–Leibler  
divergence

$$D_{KL}(P||Q) = \sum_y P(y) \log \frac{P(y)}{Q(y)}$$

$$\begin{matrix} 1.00 \\ 0.00 \\ 0.00 \end{matrix}$$

Correct  
probs

## Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Softmax  
Function

Maximize probability of correct class

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

Putting it all together:

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

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

## Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Softmax  
Function

Maximize probability of correct class

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

Putting it all together:

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

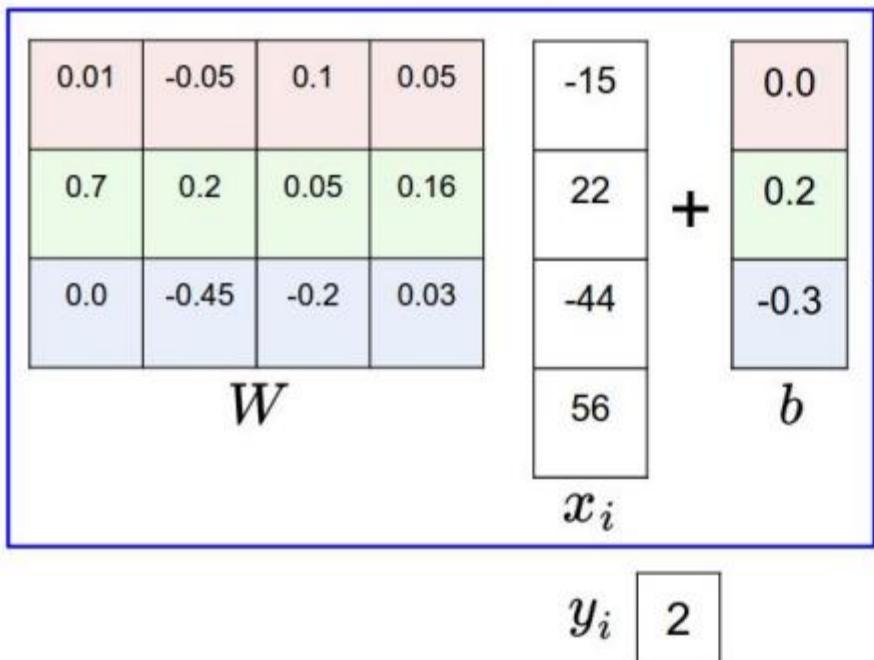
cat	<b>3.2</b>
car	5.1
frog	-1.7

Q1: What is the min/max possible softmax loss  $L_i$ ?

Q2: At initialization all  $s_j$  will be approximately equal;  
what is the softmax loss  $L_i$ , assuming  $C$  classes?

## Softmax vs. SVM

matrix multiply + bias offset



hinge loss (SVM)

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

cross-entropy loss (Softmax)

$$\begin{matrix} -2.85 \\ 0.86 \\ 0.28 \end{matrix} \xrightarrow{\text{exp}} \begin{matrix} 0.058 \\ 2.36 \\ 1.32 \end{matrix} \xrightarrow[\text{(to sum to one)}]{\text{normalize}} \begin{matrix} 0.016 \\ 0.631 \\ 0.353 \end{matrix}$$

$$-\log(0.353) = 0.452$$

### 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  $y_i = 0$

Q: What is the **softmax loss** and the **SVM** loss?

# Recap

How do we find the best  $W$ ?

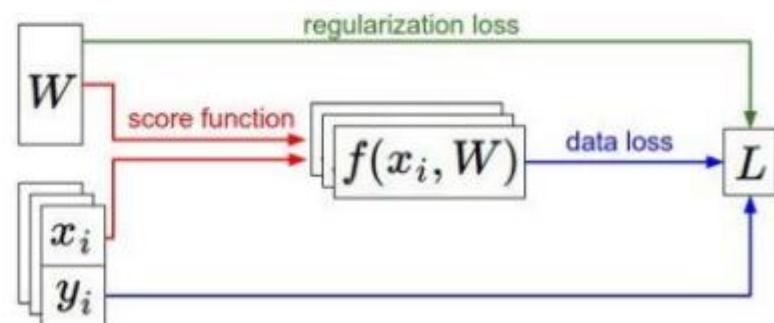
- We have some dataset of  $(x, y)$
- We have a **score function**:  $s = f(x; W) = Wx$  e.g.
- We have a **loss function**:

$$L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{sj}}\right) \quad \text{Softmax}$$

SVM

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

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



# Оптимизация параметров



## 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
```

15.5% accuracy! not bad!  
(SOTA is ~99.7%)

### Strategy #2: Follow the slope



### Strategy #2: Follow the slope

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

$\mathbf{dx}$  has same shape as  $\mathbf{x}$

$$\frac{dL(x, w)}{dw} = \lim_{h \rightarrow 0} \frac{df(x, w + h) - f(x, w)}{dh}$$

$\mathbf{dW}$  has same shape as  $\mathbf{W}$

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

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

## Оптимизация параметров

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

$\frac{dL}{dx} = \text{gradient dW:}$

[?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

**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**

$\frac{dL}{dx} = \text{gradient dW:}$

[?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

**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**

$$\frac{dL}{dx} = \text{gradient } dW:$$

**[-2.5,**  
**?,**  
**?,**

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

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

**?,**  
**?,...]**

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (second dim):**

[0.34,  
-1.11 + **0.0001**,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25353**

**gradient dW:**

[-2.5,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**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 1.25347**

**gradient dW:**

**[-2.5,  
0.6,  
?,  
?]**

$$\begin{aligned} & (1.25353 - 1.25347) / 0.0001 \\ & = 0.6 \end{aligned}$$

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

**?,...]**

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**W + h (third dim):**

[0.34,  
-1.11,  
0.78 + **0.0001**,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347****gradient dW:**

[-2.5,  
0.6,  
**0**,  
?,  
?]

**Numeric Gradient**

- Slow! Need to loop over all dimensions
- Approximate

?,...]

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

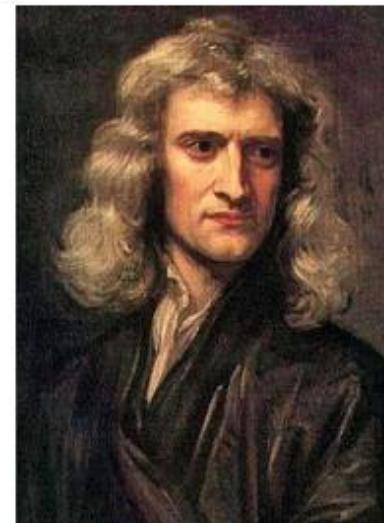
$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2$$

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

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

want  $\nabla_W L$

Use calculus to compute an  
**analytic gradient**



[This image](#) is in the public domain



[This image](#) is in the public domain

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

$dW = \dots$   
(some function  
data and W)

**gradient dW:**

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

### In summary:

- Numerical gradient: approximate, slow, easy to write
- Analytic gradient: exact, fast, error-prone

=>

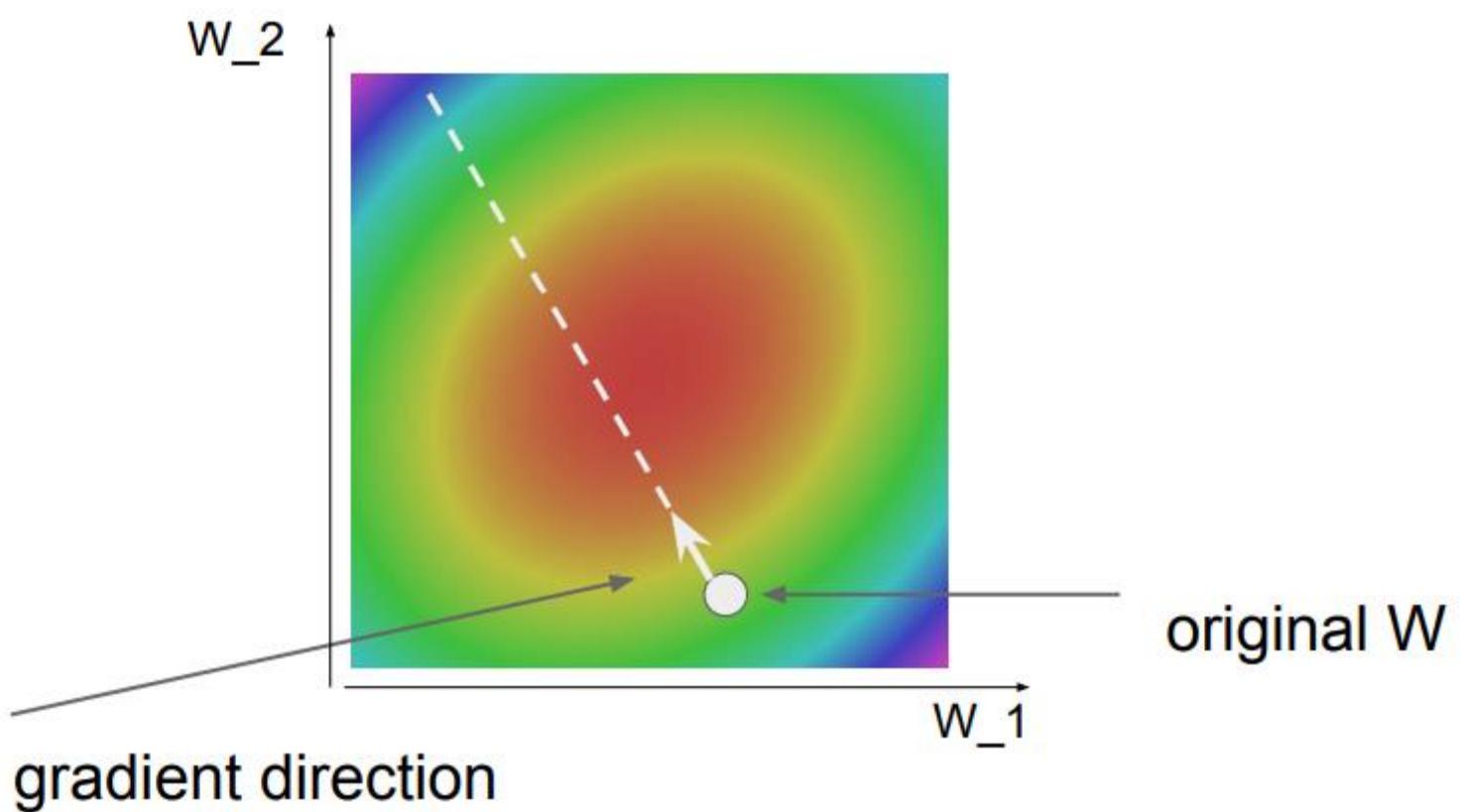
In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check**.

# Gradient Descent

```
# Vanilla Gradient Descent

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

## Оптимизация параметров



- 1. Линейные классификаторы позволяют с помощью скалярного произведения вычислить класс, к которому принадлежит изображение (вектор).**
- 2. SVM классифицирует изображения на основе его скалярного произведения с опорными векторами. Hinge loss – функция от значений этих проекций.**
- 3. Softmax классифицирует изображения на основе оценок вероятностей принадлежности к тому или иному классу. Cross entropy loss – функция от распределения вероятностей.**
- 4. Оптимальные параметры  $W$  независимо от вида loss функции вычисляются с помощью стохастического градиентного спуска.**