# Lecture 2:
# Image Classification with Linear Classifiers

Dr. José Ramón Iglesias

DSP-ASIC BUILDER GROUP
Director Semillero TRIAC
Ingenieria Electronica
Universidad Popular del Cesar

Universidad Popular del Cesar

# Administrative: Assignment 1

- K-Nearest Neighbor

- Linear classifiers: SVM, Softmax

- Two-layer neural network

- Image features

# Syllabus

| Deep Learning Basics | Convolutional Neural Networks | Computer Vision Applications |
|---|---|---|
| **Data-driven approaches** | Convolutions | RNNs / Attention / Transformers |
| **Linear classification & kNN** | PyTorch / TensorFlow | Image captioning |
| Loss functions | Activation functions | Object detection and segmentation |
| Optimization | Batch normalization | Style transfer |
| Backpropagation | Transfer learning | Video understanding |
| Multi-layer perceptrons | Data augmentation | Generative models |
| Neural Networks | Momentum / RMSProp / Adam | Self-supervised learning |
| | Architecture design | 3D vision |
| | | Human-centered AI |
| | | Fairness & ethics |

Image Classification with Linear Classifiers

Universidad Popular del Cesar

# Image Classification

A Core Task in Computer Vision

Today:
- The image classification task
- Two basic data-driven approaches to image classification
  - K-nearest neighbor and linear classifier

Universidad
Popular del Cesar

# **Image Classification**: A core task in Computer Vision

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

⟶                    cat

Image Classification with Linear Classifiers
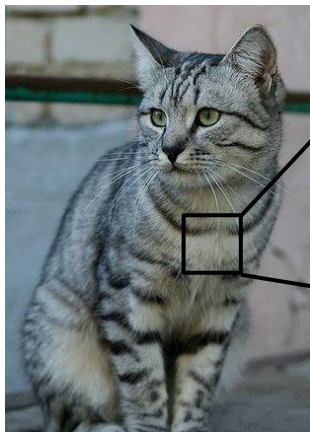
# **The Problem**: Semantic Gap
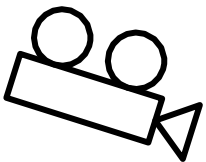


```
[[105 112 108 111 104  99 106  99  96 103 112 119 104  97  93  87]
 [ 91  98 102 106 104  79  98 103  99 105 123 136 110 105  94  85]
 [ 76  85  90 105 128 105  87  96  95  99 115 112 106 103  99  85]
 [ 99  81  81  93 120 131 127 100  95  98 102  99  96  93 101  94]
 [106  91  61  64  69  91  88  85 101 107 109  98  75  84  96  95]
 [114 108  85  55  55  69  64  54  64  87 112 129  98  74  84  91]
 [133 137 147 103  65  81  80  65  52  54  74  84 102  93  85  82]
 [128 137 144 140 109  95  86  70  62  65  63  63  60  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]
 [115 114 109 123 150 148 131 118 113 109 100  92  74  65  72  78]
 [ 89  93  90  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  80 101 124 126 119 101 107 114 131 119]
 [ 63  65  75  88  89  71  62  81 120 138 135 105  81  98 110 118]
 [ 87  65  71  87 106  95  69  45  76 130 126 107  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  83 112 153 149 122 109 104  75  80 107 112  99]
 [122 121 102  80  82  86  94 117 145 148 153 102  58  78  92 107]
 [122 164 148 103  71  56  78  83  93 103 119 139 102  61  69  84]]
```
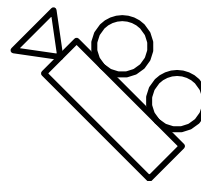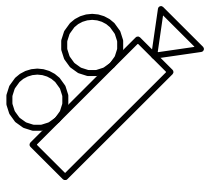
What the computer sees

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

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

Image Classification with Linear Classifiers

# **Challenges**: Viewpoint variation



All pixels change when the camera moves!

Image Classification with Linear Classifiers

# **Challenges**: Illumination

Image Classification with Linear Classifiers

**Universidad Popular del Cesar**

# **Challenges**: Background Clutter


This image is CC0 1.0 public domain


This image is CC0 1.0 public domain

Image Classification with Linear Classifiers

# **Challenges**: Occlusion

Image Classification with Linear Classifiers

# **Challenges**: Deformation

Image Classification with Linear Classifiers

Universidad Popular del Cesar

# **Challenges**: Intraclass variation

Image Classification with Linear Classifiers

# **Challenges**: Context

Image Classification with Linear Classifiers

# Modern computer vision algorithms



This image is CC0 1.0 public domain

Image Classification with Linear Classifiers

# An image classifier

```python
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

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

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

**Example training set**



airplane
automobile
bird
cat
deer

Universidad
Popular del Cesar

# Nearest Neighbor Classifier

# First classifier: **Nearest Neighbor**

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

Memorize all
data and labels

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

Predict the label
of the most similar
training image

Universidad
Popular del Cesar

# First classifier: **Nearest Neighbor**

?

deer　　bird　　plane　　cat　　car

Training data with labels

query data

Distance Metric ┆ , ┆ $\rightarrow \mathbb{R}$

# Distance Metric to compare images

**L1 distance:** $d_1(I_1, I_2) = \sum_P |I_1^p - I_2^p|$

| test image | | | |
|---|---|---|---|
| 56 | 32 | 10 | 18 |
| 90 | 23 | 128 | 133 |
| 24 | 26 | 178 | 200 |
| 2 | 0 | 255 | 220 |

| training image | | | |
|---|---|---|---|
| 10 | 20 | 24 | 17 |
| 8 | 10 | 89 | 100 |
| 12 | 16 | 178 | 170 |
| 4 | 32 | 233 | 112 |

| pixel-wise absolute value differences | | | |
|---|---|---|---|
| 46 | 12 | 14 | 1 |
| 82 | 13 | 39 | 33 |
| 12 | 10 | 0 | 30 |
| 2 | 32 | 22 | 108 |

add → 456

Image Classification with Linear Classifiers

**Universidad Popular del Cesar**

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
```

Image Classification with Linear Classifiers

Universidad Popular del Cesar

```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

Memorize training data

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

```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

For each test image:
Find closest train image
Predict label of nearest image

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

```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

**Q:** With N examples, how fast are training and prediction?

**Ans**: Train O(1), predict O(N)

This is bad: we want classifiers that are **fast** at prediction; **slow** for training is ok

Image Classification with Linear Classifiers

Universidad Popular del Cesar

```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

Many methods exist for fast / approximate nearest neighbor (beyond the scope of 231N!)

A good implementation:
https://github.com/facebookresearch/faiss

Johnson et al, "Billion-scale similarity search with GPUs", arXiv 2017

# Image Classification with Linear Classifiers

Universidad Popular del Cesar

# What does this look like?



1-nearest neighbor

# K-Nearest Neighbors

Instead of copying label from nearest neighbor,
take **majority vote** from K closest points



K = 1            K = 3            K = 5

# K-Nearest Neighbors: Distance Metric

L1 (Manhattan) distance

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

L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p \left( I_1^p - I_2^p \right)^2}$$

Universidad Popular del Cesar

# K-Nearest Neighbors: Distance Metric

L1 (Manhattan) distance

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



K = 1

L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p \left( I_1^p - I_2^p \right)^2}$$



K = 1

Image Classification with Linear Classifiers

Universidad Popular del Cesar

# K-Nearest Neighbors: try it yourself!



http://vision.stanford.edu/teaching/cs231n-demos/knn/

Image Classification with Linear Classifiers

# Hyperparameters

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

These are **hyperparameters**: choices about the algorithms themselves.

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

# Setting Hyperparameters

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

| train |
|:---:|

# Setting Hyperparameters

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

**BAD**: K = 1 always works perfectly on training data

| train |
|:---:|

# Setting Hyperparameters

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

**BAD**: K = 1 always works perfectly on training data

| train |
|:---:|

**Idea #2**: choose hyperparameters that work best on **test** data

| train | test |
|:---:|:---:|

# Setting Hyperparameters

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

**BAD**: K = 1 always works perfectly on training data

| train |
|---|

**Idea #2**: choose hyperparameters that work best on **test** data

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

| train | test |
|---|---|

Never do this!

# Setting Hyperparameters

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

**BAD**: K = 1 always works perfectly on training data

| train |
|:-----:|

**Idea #2**: choose hyperparameters that work best on **test** data

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

| train | test |
|:-----:|:----:|

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

**Better!**

| train | validation | test |
|:-----:|:----------:|:----:|

# Setting Hyperparameters

| train |
|-------|

**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 not used too frequently in deep learning

# Example Dataset: **CIFAR10**

**10** classes
**50,000** training images
**10,000** testing images

| | |
|---|---|
| **airplane** |  |
| **automobile** | |
| **bird** | |
| **cat** | |
| **deer** | |
| **dog** | |
| **frog** | |
| **horse** | |
| **ship** | |
| **truck** | |

Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images", Technical Report, 2009.

Image Classification with Linear Classifiers

# Example Dataset: **CIFAR10**

**10** classes
**50,000** training images
**10,000** testing images

Test images and nearest neighbors

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck



Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images", Technical Report, 2009.

Image Classification with Linear Classifiers

# Setting Hyperparameters



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)

# What does this look like?

# What does this look like?

Universidad
Popular del Cesar

# k-Nearest Neighbor with pixel distance **never used.**

- Distance metrics on pixels are not informative

| Original | Occluded | Shifted (1 pixel) | Tinted |



(All three images on the right have the same pixel distances to the one on the left)

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

# k-Nearest Neighbor with pixel distance **never used.**

- Curse of dimensionality

Dimensions = 1
Points = 4

Dimensions = 2
Points = $4^2$

Dimensions = 3
Points = $4^3$

# K-Nearest Neighbors: Summary

In **image classification** we start with a **training set** of images and labels, and must predict labels on the **test set**

The **K-Nearest Neighbors** classifier predicts labels based on the K nearest training examples

Distance metric and K are **hyperparameters**

Choose hyperparameters using the **validation set**;

Only run on the test set once at the very end!

Universidad
Popular del Cesar

# Linear Classifier

# Parametric Approach



Image

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

f(**x**,**W**)

**W**
parameters
or weights

**10** numbers giving
class scores

# Parametric Approach: Linear Classifier

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

**Image**



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

$f(\textcolor{blue}{\mathbf{x}},\textcolor{red}{\mathbf{W}})$

$\textcolor{red}{W}$

**parameters
or weights**

**10** numbers giving
class scores

# Parametric Approach: Linear Classifier

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

**3072x1**

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

**Image**

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

**10** numbers giving class scores

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

**W**

parameters
or weights

# Parametric Approach: Linear Classifier

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

3072x1

10x1

10x1   10x3072

## Image

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

f(**x**,**W**)

**10** numbers giving
class scores

W

parameters
or weights

Neural Network

Linear classifiers

This image is CC0 1.0 public domain

Image Classification with Linear Classifiers

[Krizhevsky et al. 2012]

Linear layers

[He et al. 2015]

Image Classification with Linear Classifiers

# Recall CIFAR10

airplane

automobile

bird

cat

deer

dog

frog

horse

ship

truck

**50,000** training images
each image is **32x32x3**

**10,000** test images.

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

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

Flatten tensors into a vector

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

Input image

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

Flatten tensors into a vector



Input image

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

+

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

b

=

| -96.8 | Cat score |
|-------|-----------|
| 437.9 | Dog score |
| 61.95 | Ship score |

Image Classification with Linear Classifiers

Universidad Popular del Cesar

# Interpreting a Linear Classifier

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

Input image

| 56 | 231 |
| 24 | 2 |

W

| 0.2 | -0.5 |
| 0.1 | 2.0 |

| 1.5 | 1.3 |
| 2.1 | 0.0 |

| 0 | .25 |
| 0.2 | -0.3 |

b

| 1.1 |

| 3.2 |

| -1.2 |

Score

| -96.8 |

| 437.9 |

| 61.95 |

# Interpreting a Linear Classifier: <u>Visual Viewpoint</u>

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

Input image

| 56 | 231 |
| 24 | 2 |

W

| 0.2 | -0.5 |
| 0.1 | 2.0 |

| 1.5 | 1.3 |
| 2.1 | 0.0 |

| 0 | .25 |
| 0.2 | -0.3 |

b

| 1.1 |

| 3.2 |

| -1.2 |

Score

| -96.8 |

| 437.9 |

| 61.95 |

plane | car | bird | cat | deer | dog | frog | horse | ship | truck

Universidad Popular del Cesar

# Interpreting a Linear Classifier: <u>Geometric Viewpoint</u>

**Analogía de las imágenes como puntos de alta dimensión.** Dado que las imágenes se estiran en vectores de columna de alta dimensión, podemos interpretar cada imagen como un solo punto en este espacio (por ejemplo, cada imagen en CIFAR-10 es un punto en un espacio de 3072 dimensiones de 32x32x3 píxeles). Análogamente, todo el conjunto de datos es un conjunto (etiquetado) de puntos.

Dado que definimos la puntuación de cada clase como una suma ponderada de todos los píxeles de la imagen, cada puntuación de clase es una función lineal sobre este espacio. No podemos visualizar espacios de 3072 dimensiones, pero si imaginamos aplastar todas esas dimensiones en solo dos dimensiones, entonces podemos tratar de visualizar lo que el clasificador podría estar haciendo:

# Interpreting a Linear Classifier: <u>Geometric Viewpoint</u>

Representación de dibujos animados del espacio de la imagen, donde cada imagen es un solo punto, y se visualizan tres clasificadores. Usando el ejemplo del clasificador de automóviles (en rojo), la línea roja muestra todos los puntos en el espacio que obtienen una puntuación de cero para la clase de automóvil. La flecha roja muestra la dirección del aumento, por lo que todos los puntos a la derecha de la línea roja tienen puntuaciones positivas (y linealmente crecientes), y todos los puntos a la izquierda tienen puntuaciones negativas (y linealmente decrecientes).

Universidad
Popular del Cesar

# Interpreting a Linear Classifier: <u>Geometric Viewpoint</u>

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

Cada fila de W es un clasificador para una de las clases. A medida que cambiamos una fila la línea que corresponde a ese pixel girará en diferentes direcciones

Los sesgos b permite que el clasificador traduzca la linea

car classifier

airplane classifier

deer classifier

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

Image Classification with Linear Classifiers

# 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

# Linear Classifier – Choose a good 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 |

TODO:

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

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

Image Classification with Linear Classifiers

**Universidad Popular del Cesar**

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** |

Universidad
Popular del Cesar

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

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

|  |  |  |  |
|------|------|------|------|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

Universidad
Popular del Cesar

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



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

|  |  |  |  |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

Universidad
Popular del Cesar

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

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



|  |  |  |  |
|------|------|------|------|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

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

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

Image Classification with Linear Classifiers

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

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: $s = f(x_i, W)$

|       |       |       |       |
|-------|-------|-------|-------|
| cat   | **3.2** | 1.3   | 2.2   |
| car   | 5.1   | **4.9** | 2.5   |
| frog  | -1.7  | 2.0   | **-3.1** |

the SVM loss has the form:

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

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

Image Classification with Linear Classifiers

**UPC** Universidad **Popular del Cesar**

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** |

Loss

$s_{y_i} - s_j$

difference in scores between correct and incorrect class

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

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

Image Classification with Linear Classifiers

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



**Interpreting Multiclass SVM loss:**



Loss

$s_{y_i} - s_j$

difference in scores between correct and incorrect class

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** |

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

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

Image Classification with Linear Classifiers

Universidad Popular del Cesar

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** |

Loss

$s_{y_i} - s_j$

difference in scores between correct and incorrect class

1

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

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

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: $s = f(x_i, W)$

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: $s = f(x_i, W)$



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

Image Classification with Linear Classifiers

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: $s = f(x_i, W)$



| | 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 | |

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

Image Classification with Linear Classifiers


Universidad Popular del Cesar

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: $s = f(x_i, W)$



| | 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 |

the SVM loss has the form:

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

= 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

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

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: $s = f(x_i, W)$



|  | 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 |

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:

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

L = (2.9 + 0 + 12.9)/3
= **5.27**

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

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

**Multiclass SVM loss:**

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



| | |
|---|---|
| cat | 1.3 |
| car | **4.9** |
| frog | 2.0 |
| Losses: | 0 |

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 ≈ 0. What is the loss $L_i$, assuming N examples and C classes?

Image Classification with Linear Classifiers

Universidad Popular del Cesar

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: $s = f(x_i, W)$

the SVM loss has the form:

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

Q4: What if the sum was over all classes? (including j = y_i)

Universidad Popular del Cesar

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

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: $s = f(x_i, W)$

|  | | | |
|---|---|---|---|
| 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 |

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?

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

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

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: $s = f(x_i, W)$



|       |      |      |      |
|-------|------|------|------|
| 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 |

the SVM loss has the form:

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

Q6: What if we used

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

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

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

**Multiclass SVM loss:**



Loss

$s_{y_i} - s_j$

difference in scores between correct and incorrect class

1

|        |      |     |      |
|--------|------|-----|------|
| 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 |

Q6: What if we used

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

Image Classification with Linear Classifiers

Universidad
Popular del Cesar

# Multiclass SVM Loss: Example 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
```

# First calculate scores
# Then calculate the margins $s_j$ - $s_{yi}$ + 1
# only sum j is not $y_i$, so when j = $y_i$, set to zero.
# sum across all j

**Universidad Popular del Cesar**

# Softmax classifier

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



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

cat **3.2**

car 5.1

frog -1.7

Universidad
Popular del Cesar

# **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     **3.2**

car     5.1

frog    -1.7

**Universidad Popular del Cesar**

# **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 >= 0

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

unnormalized 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 >= 0

Probabilities must sum to 1

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

exp → normalize →

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 >= 0

Probabilities must sum to 1

|  | Unnormalized log-probabilities / logits | | unnormalized probabilities | | probabilities |
|---|---|---|---|---|---|
| cat | **3.2** | | **24.5** | | **0.13** |
| car | 5.1 | exp | 164.0 | normalize | 0.87 |
| frog | -1.7 | | 0.18 | | 0.00 |

Image Classification with Linear Classifiers

**Universidad Popular del Cesar**

# **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 >= 0

Probabilities must sum to 1

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

|  | Unnormalized log-probabilities / logits | | unnormalized probabilities | | probabilities | |
|---|---|---|---|---|---|---|
| cat | **3.2** | | **24.5** | | **0.13** | $L_i = -\log(0.13)$ = **2.04** |
| car | 5.1 | exp | 164.0 | normalize | 0.87 | |
| frog | -1.7 | | 0.18 | | 0.00 | |

**Universidad Popular del Cesar**

# **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 >= 0

Probabilities must sum to 1

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

cat

| | | |
|---|---|---|
| **3.2** | **24.5** | **0.13** |
| 5.1 | 164.0 | 0.87 |
| -1.7 | 0.18 | 0.00 |

exp

normalize

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

cat

car

frog

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

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

**Universidad Popular del Cesar**

# **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

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

Probabilities must be >= 0

Probabilities must sum to 1

|  | logits | $\xrightarrow{\text{exp}}$ | unnormalized probabilities | $\xrightarrow{\text{normalize}}$ | probabilities | compare | Correct probs |
|---|---|---|---|---|---|---|---|
| cat | **3.2** | | **24.5** | | **0.13** | | **1.00** |
| car | 5.1 | | 164.0 | | 0.87 | | 0.00 |
| frog | -1.7 | | 0.18 | | 0.00 | | 0.00 |

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

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

Probabilities must be >= 0

Probabilities must sum to 1

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

|  | logits | exp | probabilities | compare | Correct probs |
|---|---|---|---|---|---|
| cat | **3.2** | **24.5** | **0.13** | | **1.00** |
| car | 5.1 | 164.0 | 0.87 | | 0.00 |
| frog | -1.7 | 0.18 | 0.00 | | 0.00 |

exp → normalize → compare

*Kullback–Leibler divergence*

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

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

Correct probs

Image Classification with Linear Classifiers

Universidad Popular del Cesar

March 31, 2022

# **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

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

Probabilities must be >= 0

Probabilities must sum to 1

|      | logits | exp   | normalize | compare | Correct probs |
|------|--------|-------|-----------|---------|---------------|
| cat  | **3.2** | **24.5** | **0.13** |         | **1.00** |
| car  | 5.1    | 164.0 | 0.87      |         | 0.00 |
| frog | -1.7   | 0.18  | 0.00      |         | 0.00 |

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

*Cross Entropy*

$$H(P, Q) =$$
$$H(p) + D_{KL}(P \| Q)$$

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      **3.2**

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 **3.2**

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 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

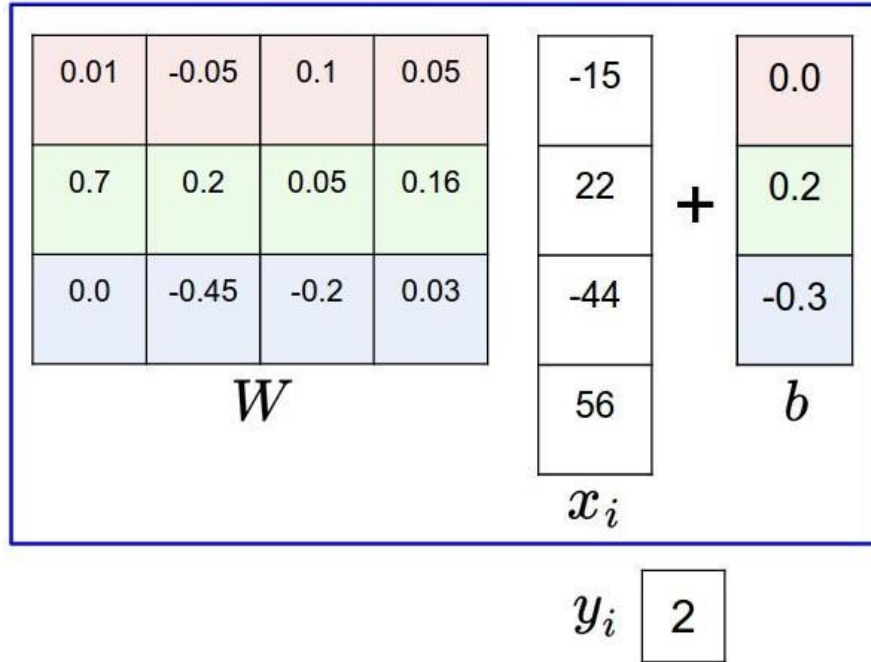$$L_i = -\log P(Y = y_i | X = x_i)$$

Putting it all together:

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

cat **3.2**

car 5.1

frog -1.7

Q2: At initialization all s will be approximately equal; what is the loss?
A: -log(1/C) = log(C),
If C = 10, then $L_i$ = log(10) ≈ 2.3

# Softmax vs. SVM

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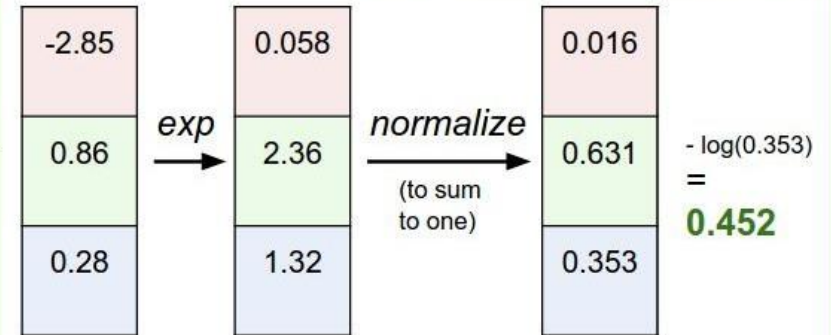
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

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

**Universidad Popular del Cesar**

# 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?

**Universidad Popular del Cesar**

# 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:
[20, -2, 3]
[20, 9, 9]
[20, -100, -100]
and $\quad y_i = 0$

Q: What is the **softmax loss** and the **SVM** loss **if I double the correct class score from 10 -> 20**?

**Universidad Popular del Cesar**

# Coming up:

- Regularization
- Optimization

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

Universidad Popular del Cesar