# CENG 506 Deep Learning

Lecture 2 – Linear Classification
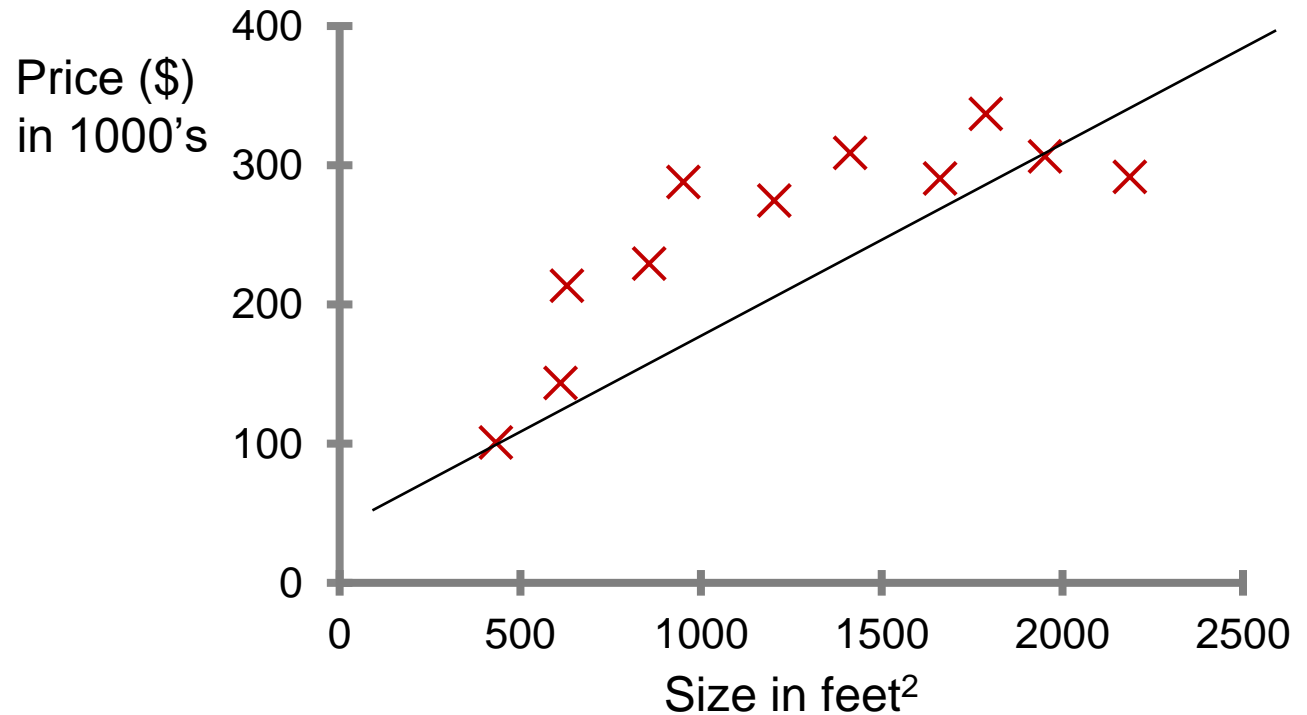
(Logistic Regression, Multi-class Classification,

Softmax Loss, Regularization, Learning Rate)
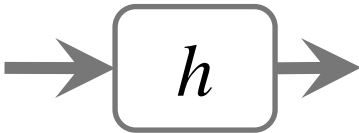
# Recap

Last week we talked about linear regression problem.
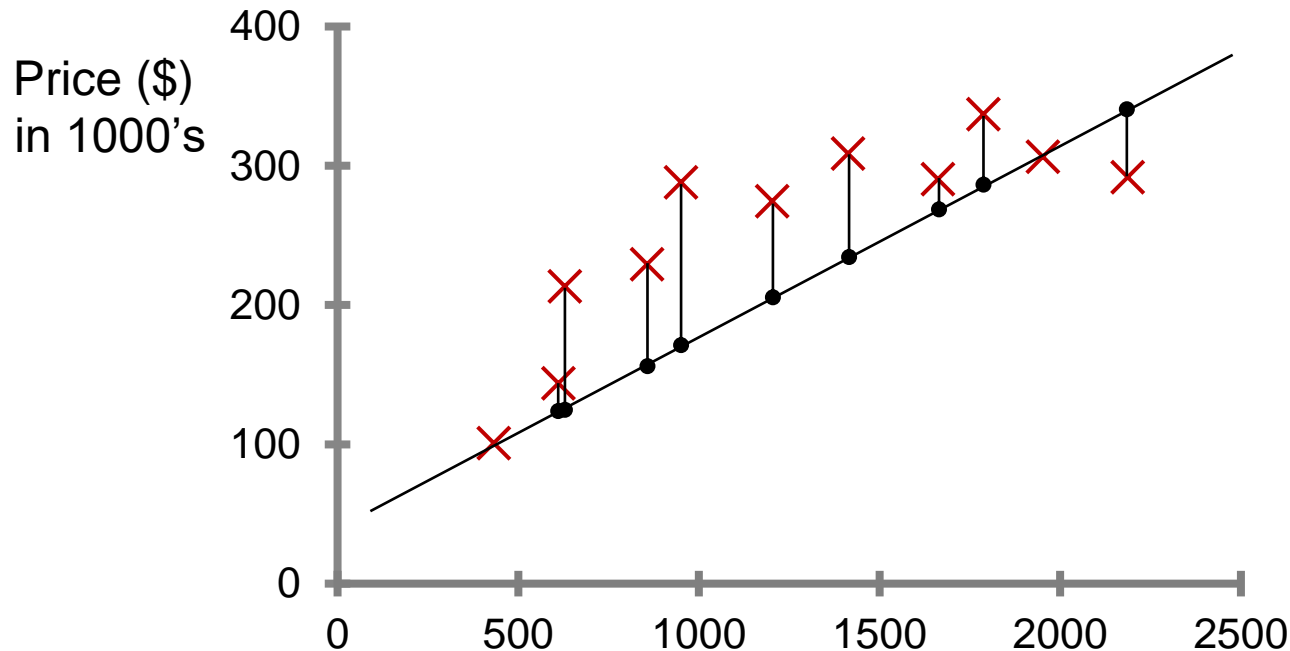
Example:

House prices according to house size

Price ($) in 1000's



Size in feet$^2$

Size of house $\rightarrow$ $h$ $\rightarrow$ Estimated price

$$h_w(x) = w_0 + w_1 x$$

# Recap

We defined a cost function to minimize.



Price ($)
in 1000's

$y$: real prices
(red crosses)

Our hypothesis:
$h(x) = w_0 + w_1 x$

First sample cost:
$(y^{(1)} - h(x^{(1)}))^2$

Total cost:
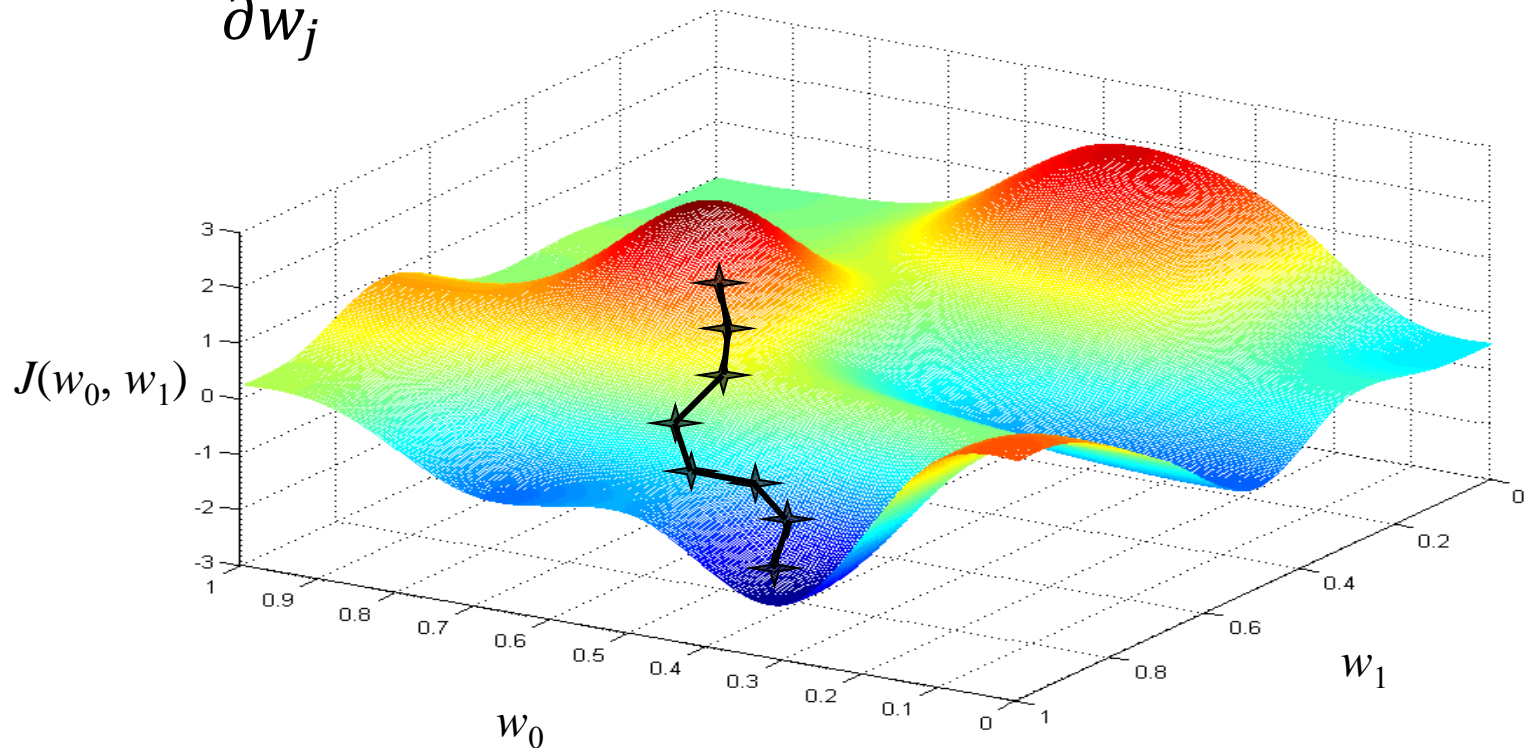$\Sigma (y^{(i)} - h(x^{(i)}))^2$

# Recap

We changed $w_0, w_1$ with gradient descent to reach min cost.

Repeat until convergence {

$$w_j := w_j - \alpha \frac{\partial J(w_0, w_1)}{\partial w_j}$$
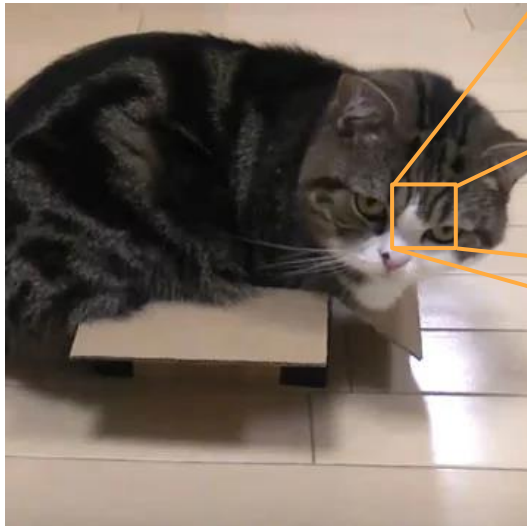
}

# Binary Classification



128x128x3 image

$$y \in \{0,1\}$$

0: Non-cat
1: Cat

# Binary Classification



128x128x3 = 49152. Data vector becomes

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{49152} \end{bmatrix} \qquad y = 1$$

# Logistic Regression

$$h_w(x) = g(w^T \cdot x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$

$$h_w(x) = \frac{1}{1+e^{-w^T \cdot x}}$$

$g(z)$ is called
*Logistic function*
or
*Sigmoid fuction*



$$h_w(x) = g(w^T \cdot x)$$

$z = w^T \cdot x < 0$  $z = w^T \cdot x > 0$

Parameter vector $w$ has the same size with $x$.
Parameters are learned (fitted) with Log. Reg. algorithm.

# Decision Boundary

One strategy for decision is:

Predict $y$=1 if $h_\theta(x) > 0.5$

   $y$=0 if $h_\theta(x) <= 0.5$

This means

Predict $y$=1 if $\theta^T \cdot x > 0$

   $y$=0 if $\theta^T \cdot x <= 0$

$$h_\theta(x) = g(\theta^T \cdot x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$



Decision
boundary

# Logistic Regression

For an image (e.g. 49152 dimensions) we can not visualize, but it is a linear boundary.
An example in 2D looks like this:

$$h_w(x) = g(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2)$$

$$h_w(x) = g(-3 + 1 \cdot x_1 + 1 \cdot x_2)$$



$x_1 + x_2 = 3$ : Decision boundary

# Logistic Regression

How to estimate the best $w$ for the decision boundary?

We could minimize a simple cost function as follows:

$$J(w) = \frac{1}{2m}\sum_{i=1}^{m}\left(h_w(x^{(i)}) - y^{(i)}\right)^2$$

It turns out that, because of the non-linearity of the sigmoid function that we use in $h_w(x)$, this cost function becomes non-convex (i.e. contains local minima).

"non-convex"

# Logistic Regression

We need another function which is convex.
Let's use **cross-entropy loss**:

$$\text{Cost} \left( h_w(x), y \right) = \begin{cases} -\log(h_w(x)) & \text{if } y=1 \\ -\log(1-h_w(x)) & \text{if } y=0 \end{cases}$$

Cost

If $y = 1$

Cost=0 @ $h_w(x)=1$ (since $y=1$)
But, cost $\rightarrow \infty$ as $h_w(x) \rightarrow 0$
which means: if we predit 0
when $y=1$, we penalize by a
very large cost.

0                    1    $h_w(x)$

# Logistic Regression

**Cross-entropy loss**:

$$\text{Cost} \ ( \ h_w(x), y \ ) = \begin{cases} - \log(h_w(x)) & \text{if } y\text{=1} \\ -\log(1\text{-}h_w(x)) & \text{if } y\text{=0} \end{cases}$$

If $y = 0$

Cost



0          1     $h_w(x)$

-log(1-z)



z

1

Cost=0 @ $h_w(x)$=0 (since $y$=0)
But, Cost $\rightarrow \infty$  as  $h_w(x) \rightarrow 1$

# Logistic Regression

$$\text{Cost } ( \, h_w(x), \, y \, ) = \begin{cases} -\log(h_w(x)) & \text{if } y=1 \\ \\ -\log(1\text{-}h_w(x)) & \text{if } y=0 \end{cases}$$

Knowing that $y$ is always equal to $0$ or $1$, we can define the cost function with a single line:

$$\text{Cost } ( \, h_w(x), \, y \, ) = - y \cdot \log(h_w(x)) - (1 - y) \cdot \log(1\text{-}h_w(x))$$

if $y=1$, cost becomes $-\log(h_w(x))$

if $y=0$, cost becomes $-\log(1\text{-}h_w(x))$

# Derivative of the Cost Function

$$J(w) = \frac{1}{m} \sum_{i=1}^{m} (-y^{(i)} \log(h_w(x^{(i)})) - (1-y) \log(1 - h_w(x^{(i)}))$$

To estimate parameters, $\min_w J(w)$

Derivative computed using calculus is as follows:

$$\frac{\partial J(w)}{\partial w_j} = \frac{1}{m} \sum_{i=1}^{m} (h_w(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

# Gradient Descent

Repeat $\{$

$$w_j := w_j - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_w\left(x^{(i)}\right) - y^{(i)} \right) \cdot x_j^{(i)}$$

$h_w(x) = g(w^T x) = \dfrac{1}{1 + e^{-w^T x}}$

gradient

$\}$

# Gradient Descent – Learning Rate

$$w_j := w_j - \alpha \frac{\partial J(w)}{w_j}$$

If $\alpha$ is too small, gradient descent can be slow.

If $\alpha$ is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

# Multi-class Classification

Now we have a different parameter set $w$ for each class.

$w$'s are stacked up, forms a matrix $W$.

Here, $b$ is actually $w_0$, and $W$ includes $w_1, w_2$, etc.

Let's say there are $K$ classes.

image    parameters (or weights)    bias

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

Result: $K$ x 1
i.e. $K$ numbers indicating class scores

Image size: [128x128x3]
A vector of length 49152

$K$ x 49152

49152 x 1

$K$ x 1

# Multi-class Classification

Example with a 2x2 image and 3 classes:

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

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

$+$

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

$=$

| -96.8 | Class 1 score: $s_1$ |
|-------|----------------------|
| 437.9 | Class 2 score: $s_2$ |
| 61.95 | Class 3 score: $s_3$ |

$$W \qquad x \qquad b \qquad \mathrm{f}(x, W, b)$$

# Multi-class Classification

What does a linear classifier do?



airplane classifier

car classifier

deer classifier

# Softmax Function

In binary classification, sigmoid function was enough to convert the scores into two probabilities where they add up to 1.

$$P(y{=}0/\ x,\ w\ ) + P(y{=}1/\ x,\ w\ ) = 1.$$

In multi-class classification, label $y$ can take on $K$ different values.
Let $s_j = f(x,\ w_j)$ , i.e. score for the $j^{\text{th}}$ class.
Given a test input $x$, we want our hypothesis to estimate the probabilities for each value of $j = 1,..,K$ (probability of each class). We need a $K$ dimensional vector whose elements sum up to 1.

$$\begin{bmatrix} P(y = 1|x,W) \\ P(y = 2|x,W) \\ \vdots \\ P(y = K|x,W) \end{bmatrix} = \frac{1}{\sum_{j=1}^{K} e^{s_j}} \begin{bmatrix} e^{s_1} \\ e^{s_2} \\ \vdots \\ e^{s_K} \end{bmatrix} \longleftarrow \text{Softmax function}$$

# Multi-class Cross-Entropy Loss

Remember the loss in binary classification:

$$J(w) = \frac{1}{m}\sum_{i=1}^{m}(-y^{(i)}\log(h_w(x^{(i)})) - (1-y)\log(1-h_w(x^{(i)}))$$

With softmax, it can be generalized to *K* classes:

$$J(W) = \frac{1}{m}\sum_{i=1}^{m}\sum_{k=1}^{K} -1\{y^{(i)} = k\}\log\frac{e^{s_k}}{\sum_{j=1}^{K}e^{s_j}}$$

where $1\{\cdot\}$ is the indicator function such that

$$1\{a\ true\ statement\} = 1 \quad \text{and} \quad 1\{a\ false\ statement\} = 0$$

# Softmax Loss

All together, softmax loss for sample $i$ :

$$L_i = -\log \frac{e^{s_{y^{(i)}}}}{\sum_{j=1}^{K} e^{s_j}}$$

where $y^{(i)}$ is the correct class label.

# Softmax Classifier Example

Suppose 3 classes. With some $W$, the scores f($x$,$W$,$b$) are:

$$L_i = -\log \frac{e^{s_{y(i)}}}{\sum_{j=1}^{K} e^{s_j}}$$

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

Observe the probabilistic
interpretation in softmax classifier

# Softmax Classifier Example continued

- What is the min/max possible loss?

- Usually at initialization $W$ are small numbers, so all $s \sim= 0$. What is the loss?

$$L_i = -\log \frac{e^{s_{y^{(i)}}}}{\sum_{j=1}^{K} e^{s_j}}$$

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

# Softmax Classifier Example continued



|      |       |       |     |        |     |       |      |
|------|-------|-------|-----|--------|-----|-------|------|
| cat  | **3.2** | 1.3   |     | 3.7    |     | 0.03  |      |
| car  | 5.1   | **4.9** | exp | **134.3** | norm. | **0.92** | $L_i = -\log(0.92) = 0.04$ |
| frog | -1.7  | 2.0   |     | 7.4    |     | 0.05  |      |
| loss | 0.89  |       |     |        |     |       |      |

# Softmax Classifier Example continued



| | | | | | |
|---|---|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 | 9.03 | 0.43 |
| car | 5.1 | **4.9** | 2.5 | 12.18 | 0.57 |
| frog | -1.7 | 2.0 | **-3.1** | **0.05** | **0.002** |
| loss | 0.89 | 0.04 | | | |

exp ⟶ norm. ⟶

$L_i = -\log(0.002) = 2.7$

# Regularization

- We want smaller weights, results in smoother models

regularization strength

$$J(W) = \underbrace{\frac{1}{m}\sum_{i=1}^{m}\sum_{k=1}^{K} 1\{y^{(i)} = k\}\log\frac{e^{s_k}}{\sum_{j=1}^{K}e^{s_j}}}_{\text{data loss}} + \underbrace{\lambda R(W)}_{\text{regularization loss}}$$

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

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

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

- Dropout   (used with NNs, will see later)

# Regularization

- Regularization means, the coefficients ($w$) are chosen so that not only to predict well on the training data, but also the magnitude of coefficients to be as small as possible.
- Intuitively, it means that no input dimension can have a very large influence on the scores all by itself.
- Motivation (L2 regularization example):

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

$$w_A = [1,0,0,0]$$

$$w_B = [0.25,0.25,0.25,0.25]$$

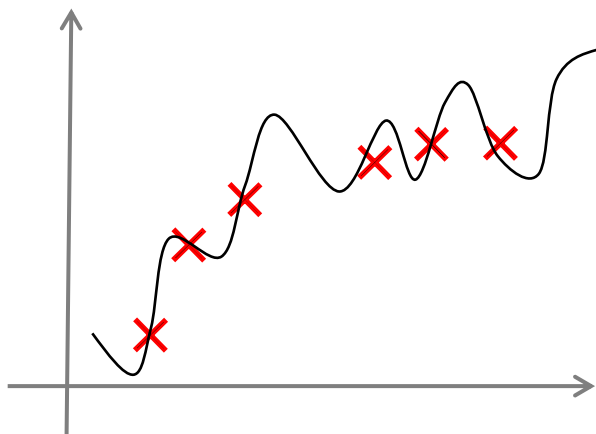$$w_A^\mathsf{T}x = w_B^\mathsf{T}x = 1$$

preferable

L2 penalty of $w_A$ is 1.0 while the L2 penalty of $w_B$ is 0.25.
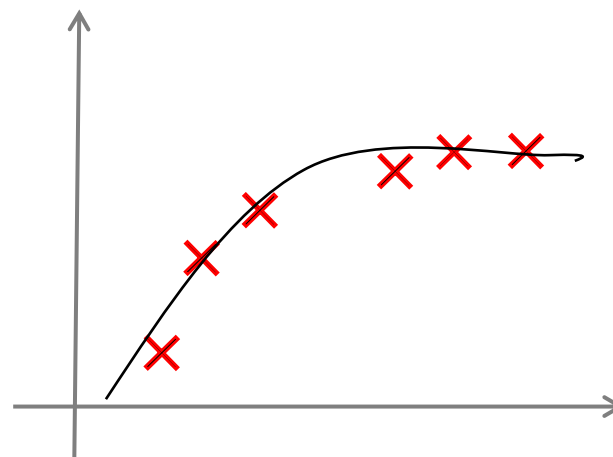
- Regularization helps a model to avoid <u>overfitting</u>.

# Regularization
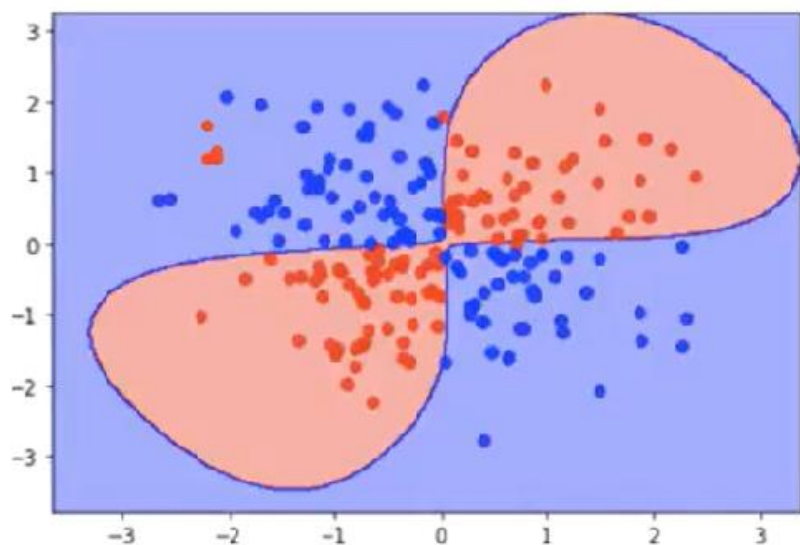
Regularization restricts a model to avoid <u>overfitting</u>.



Not regularized well.

Does not work well on test data.

Needs more regularization
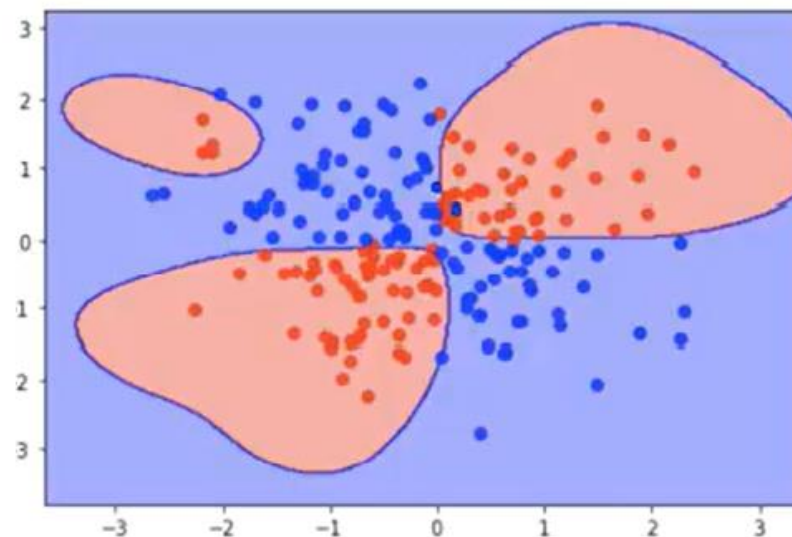
($\lambda$ to be increased)

Regularized well.

Works well on test data.

# Regularization

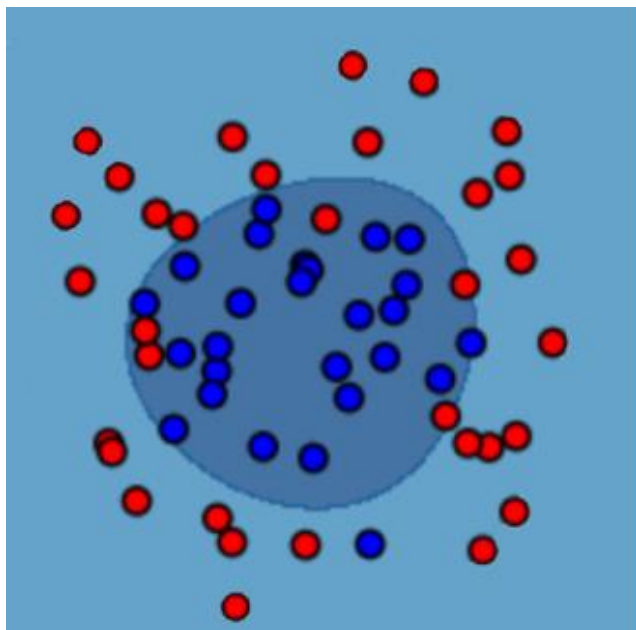Which one has the best regularization (left or right)?
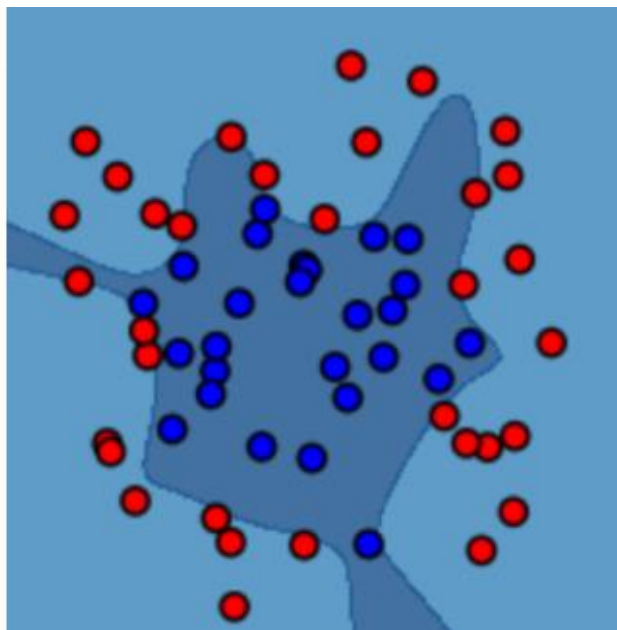


Good amount of regularization

Needs more regularization
($\lambda$ to be increased)

# Regularization

Which one has the best regularization (left or right)?



Good amount of regularization

Needs more regularization
($\lambda$ to be increased)

# Linear Classification Loss Visualization

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