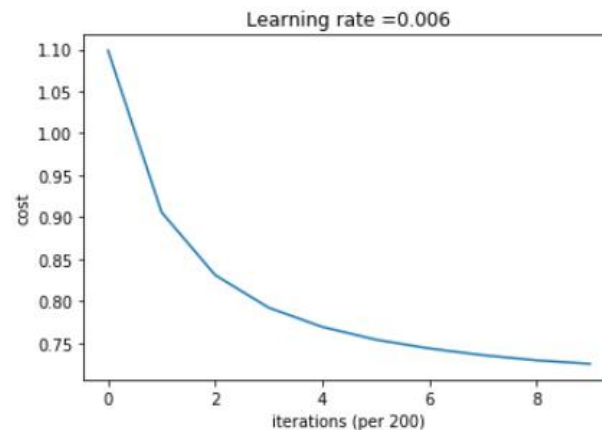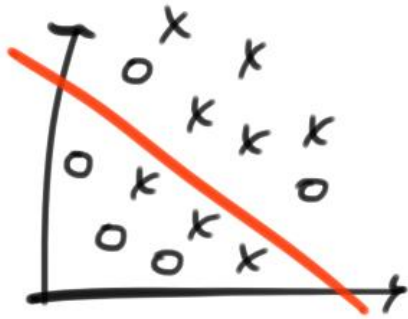- 2000  iteration

# ADAM

```
Cost after iteration 0: 1.098627
Cost after iteration 200: 0.906004
Cost after iteration 400: 0.831171
Cost after iteration 600: 0.792285
Cost after iteration 800: 0.769404
Cost after iteration 1000: 0.754191
Cost after iteration 1200: 0.743681
Cost after iteration 1400: 0.735736
Cost after iteration 1600: 0.729602
Cost after iteration 1800: 0.725414
```
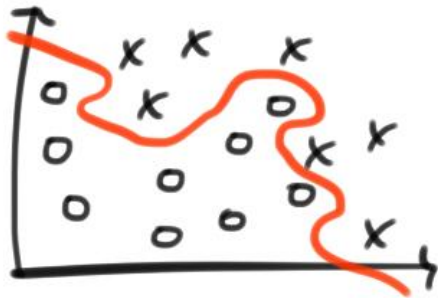
Learning rate =0.006



```
train accuracy :  0.86
test accuracy :  0.34
```
overfitting

- underfitting
  ⟺ high bias

  ⎰ ① pick network
  ⎱ ② train longer



- overfitting
  ⟺ high variance

  ⎰ ① more data
  ⎱ ② regularization

$$w = w - \alpha \, dw$$

$$J = \frac{1}{m} \sum L(\tilde{y}, y)$$

$\downarrow$ L2 regularization

$$w = w - \alpha \left( dw + \frac{\lambda}{m} w \right)$$

$$J = \frac{1}{m} \sum L(\tilde{y}, y) + \frac{\lambda}{2m} \|w\|_2^2$$

penalty

# ⊕ Adam

$$\begin{cases} v_{dW^{[l]}} = \beta_1 v_{dW^{[l]}} + (1 - \beta_1)\frac{\partial J}{\partial W^{[l]}} \\[2mm] v_{dW^{[l]}}^{corrected} = \frac{v_{dW^{[l]}}}{1-(\beta_1)^t} \\[2mm] s_{dW^{[l]}} = \beta_2 s_{dW^{[l]}} + (1 - \beta_2)(\frac{\partial J}{\partial W^{[l]}})^2 \\[2mm] s_{dW^{[l]}}^{corrected} = \frac{s_{dW^{[l]}}}{1-(\beta_2)^t} \\[2mm] W^{[l]} = W^{[l]} - \alpha\frac{v_{dW^{[l]}}^{corrected}}{\sqrt{s_{dW^{[l]}}^{corrected}} + \varepsilon} \end{cases}$$

"decoupled weight decay regularization"

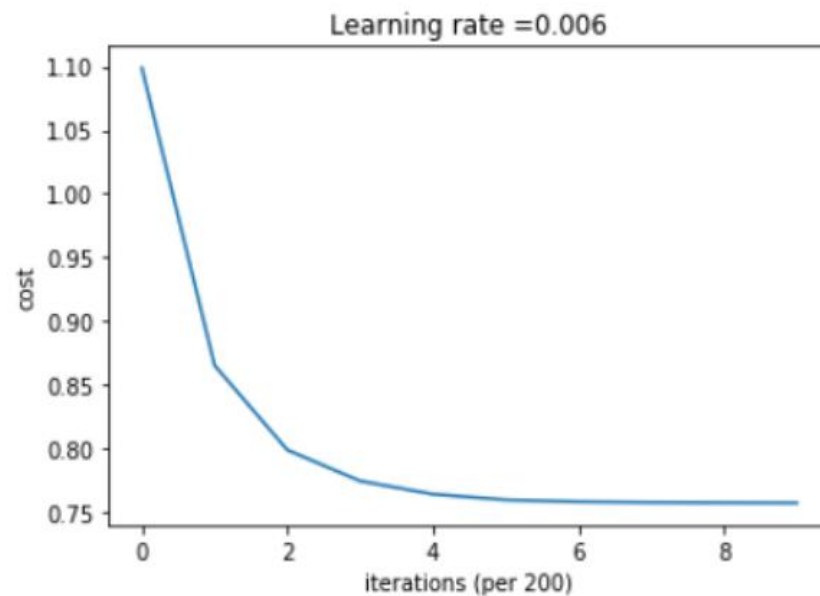① Adam은 손실함수에 L2 norm을 더하여 최적화 해도 일반화 효과를 똑바르게 된다

② weight 식에 직접 weight decay term을 추가하여 이 문제를 해결

**Algorithm 2** Adam with $L_2$ regularization and Adam with decoupled weight decay (AdamW)

1: **given** $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \lambda \in \mathbb{R}$
2: **initialize** time step $t \leftarrow 0$, parameter vector $\boldsymbol{\theta}_{t=0} \in \mathbb{R}^n$, first moment vector $\boldsymbol{m}_{t=0} \leftarrow \boldsymbol{0}$, second moment vector $\boldsymbol{v}_{t=0} \leftarrow \boldsymbol{0}$, schedule multiplier $\eta_{t=0} \in \mathbb{R}$
3: **repeat**
4:     $t \leftarrow t + 1$
5:     $\nabla f_t(\boldsymbol{\theta}_{t-1}) \leftarrow \text{SelectBatch}(\boldsymbol{\theta}_{t-1})$       ▷ select batch and return the corresponding gradient
6:     $\boldsymbol{g}_t \leftarrow \nabla f_t(\boldsymbol{\theta}_{t-1})\ +\lambda\boldsymbol{\theta}_{t-1}$
7:     $\boldsymbol{m}_t \leftarrow \beta_1 \boldsymbol{m}_{t-1} + (1 - \beta_1)\boldsymbol{g}_t$       ▷ here and below all operations are element-wise
8:     $\boldsymbol{v}_t \leftarrow \beta_2 \boldsymbol{v}_{t-1} + (1 - \beta_2)\boldsymbol{g}_t^2$
9:     $\hat{\boldsymbol{m}}_t \leftarrow \boldsymbol{m}_t/(1 - \beta_1^t)$       ▷ $\beta_1$ is taken to the power of $t$
10:    $\hat{\boldsymbol{v}}_t \leftarrow \boldsymbol{v}_t/(1 - \beta_2^t)$       ▷ $\beta_2$ is taken to the power of $t$
11:    $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$       ▷ can be fixed, decay, or also be used for warm restarts
12:    $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \eta_t \left( \alpha\hat{\boldsymbol{m}}_t/(\sqrt{\hat{\boldsymbol{v}}_t} + \epsilon)\ +\lambda\boldsymbol{\theta}_{t-1} \right)$
13: **until** *stopping criterion is met*
14: **return** optimized parameters $\boldsymbol{\theta}_t$

```
Cost after iteration 0: 1.098858
Cost after iteration 200: 0.865402
Cost after iteration 400: 0.799049
Cost after iteration 600: 0.774712
Cost after iteration 800: 0.764440
Cost after iteration 1000: 0.759916
Cost after iteration 1200: 0.758445
Cost after iteration 1400: 0.757904
Cost after iteration 1600: 0.757688
Cost after iteration 1800: 0.757495
```

```
train accuracy :   0.8583333333333333
test accuracy :   0.32
```



Learning rate =0.006

# "Double backpropagation"

- **DataGrad** (Double Backpropagation): penalize the $L2$ norm of the gradient of the original loss term with respect to the inputs.

$$L_{DG}(x, y, \Theta) = L(x, y, \Theta) + \lambda \|(\frac{\partial}{\partial x} L(x, y, \Theta))\|_2$$

```
Cost after iteration 0: 1.098838
Cost after iteration 200: 0.859740
Cost after iteration 400: 0.770849
Cost after iteration 600: 0.724785
Cost after iteration 800: 0.699078
Cost after iteration 1000: 0.681866
Cost after iteration 1200: 0.670880
Cost after iteration 1400: 0.663346
Cost after iteration 1600: 0.658710
Cost after iteration 1800: 0.655125
```

```
train accuracy :  0.85
test accuracy :  0.36
```



Learning rate =0.006