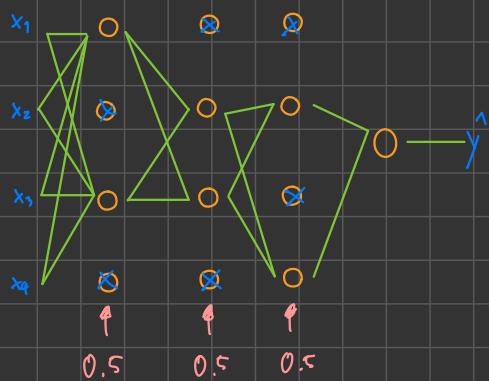
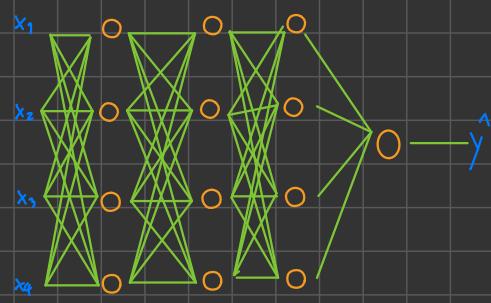


Optimize NN, #2

Dropout regularization



How to implementify dropout

- Inverted dropout

Illustrify with layer L⁴, keep_prob = 0.5

$d_3 = np.random.rand(a_3.shape[0], a_3.shape[1]) < \text{keep_prob}$

$a_3' = p * \text{multiply}(a_3, d_3)$ $\times a_3 + d_3$

$a_3' \neq \text{keep_prob} \rightarrow$ to not change expect value of next layer \rightarrow so units \rightarrow no unit shut off

$$z^{(4)} = w^{(4)} \cdot a^{(3)} + b^{(4)}$$

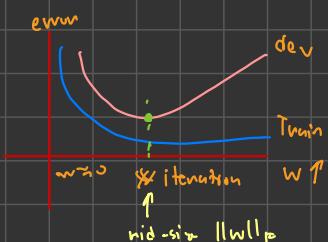
reduces by 20%
 $\therefore 0.8$

why does it work

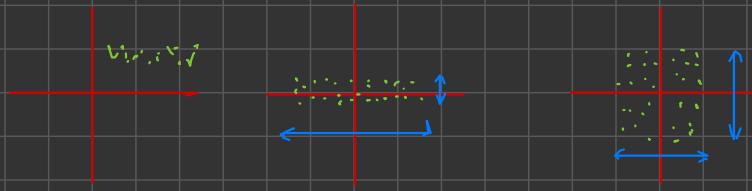
- Can't rely on any one factor, so have to spread out weights.
 - {
 - ↳ Shrink weight
- Recommend in Computer Vision \rightarrow because image have many pixels mean more iteration so dropout

Other regularization methods

- Data arguments
 - rotate
 - flip
- Early stopping
 - plot test and Dev set error when find minimum error dev set stop training



Normalizing



subtract mean:

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

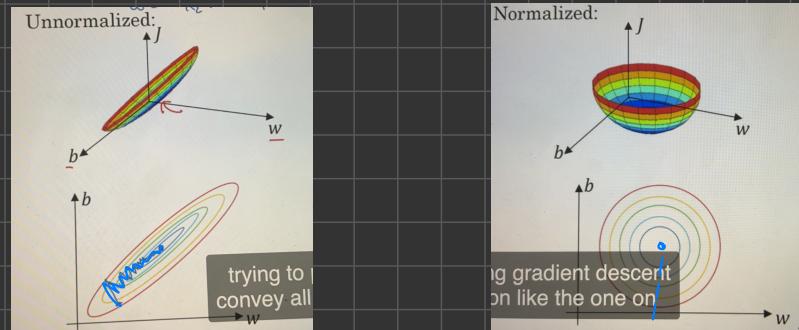
$$x = x - \mu$$

Normal variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

$$x = x / \sigma^2$$

Why we normalize



- The gradient descent can step easier (not oscillate)
- Prevent one feature have more priority than others (similar scale)
- easy to optimize ($1 + 1000$ to $0 \rightarrow 1$, $-1 \rightarrow 1$)

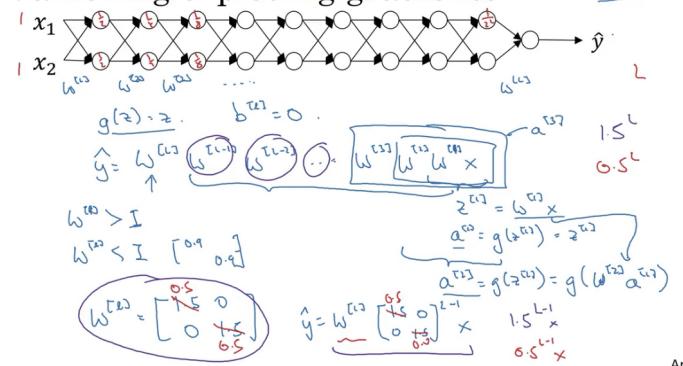
Vanishing / Exploding gradients

13:32 Fri 22 May



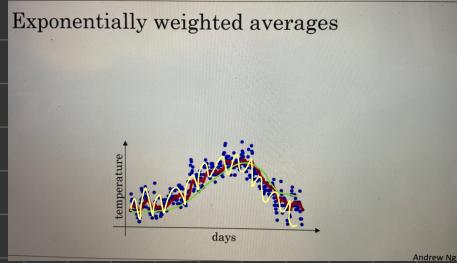
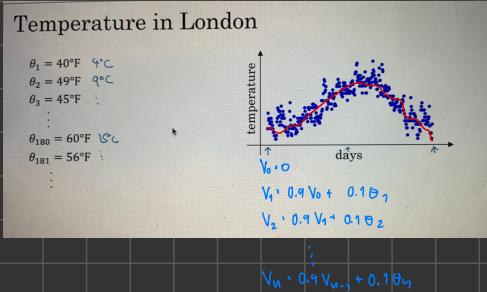
Vanishing / Exploding gradients

Vanishing/exploding gradients



Andrew Ng

Exponentially weighted averages



$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$\beta = 0.9 \approx 10 \text{ days}$$

$$\beta = 0.99 \approx 20 \text{ days}$$

$$\beta = 0.5 \approx 2 \text{ days}$$

V_t as approximately average over $\frac{1}{1-\beta}$ days temperature

$$\approx \frac{1}{1-0.99} = 50$$

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$V_{t+1} = \beta V_t + (1-\beta) \theta_{t+1}$$

$$V_{t+2} = \beta V_{t+1} + (1-\beta) \theta_{t+2}$$

$$V_{t+3} = \beta V_{t+2} + (1-\beta) \theta_{t+3}$$

$$\vdots$$

$$V_{t+100} = \beta V_{t+99} + (1-\beta) \theta_{t+100}$$

$$= \beta^1 \theta_t + \beta^2 \theta_{t+1} + \beta^3 \theta_{t+2} + \dots + \beta^{100} \theta_{t+99}$$

$$+ \dots$$

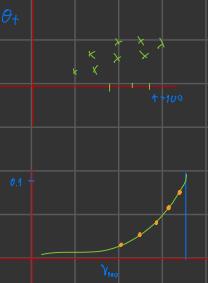
$$\beta^{100} \approx 0.0001 \quad (1-\beta)^{100} \approx 1$$

The peak to change

$$\text{When } \beta^0 \text{ will be } 1$$

$$0.99^{\frac{1}{0.0001}} \rightarrow 2.48^{\frac{1}{0.0001}} \approx 1$$

that a peak of change



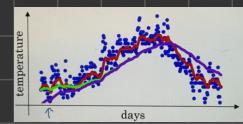
Implementing

```

V_t := 0
V_t := beta * v_t + (1-beta) * theta_t
V_t := beta * v_t + (1-beta) * theta_t
:
v_t := 0
Report & Get next theta_t
    V_t := beta * V_t + (1-beta) * theta_t
    :

```

Bias Correction



When add more to average the windows will move right and slow to change

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$V_t = 0$$

$$V_t = 0.9(0) + 0.02 \theta_1$$

$$V_t = 0.9(0.9(0) + 0.02 \theta_2) + 0.02 \theta_1$$

$$V_t = 0.9(0.9(0.9(0) + 0.02 \theta_3) + 0.02 \theta_2) + 0.02 \theta_1$$

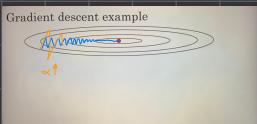
$$\vdots$$

V_t This helps are exponentially weight to help eliminate value slowly on

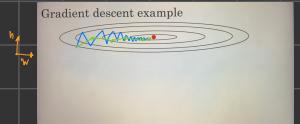
$$1 - \beta^t = 1 - (0.99)^t = 0.0001t$$

$$\frac{V_t}{1 - \beta^t} = \frac{0.99t \theta_t + 0.0001t}{0.0001t} = 99 \theta_t + 1$$

Momentum



RMSProp



Momentum

On iteration t :

Compute $\delta w, \delta b$ on current min-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) \delta w, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) \delta b$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) \delta w^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) \delta b^2$$

friction \rightarrow Velocity

$$w := w - \alpha V_{dw}, \quad b := b - \alpha V_{db}$$

Hyperparameters: α, β_1, β_2

Adam optimization

$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$

On iteration t :

Compute $\delta w, \delta b$ using current minibatch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) \delta w, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) \delta b \quad \leftarrow "momentum" \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) \delta w^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) \delta b^2 \quad \leftarrow "bias\ correction" \beta_2$$

$$V_{dw}^{current} = V_{dw} / (\beta_1 \beta_2), \quad V_{db}^{current} = V_{db} / (\beta_1 \beta_2)$$

$$S_{dw}^{current} = S_{dw} / (\beta_2 \beta_2), \quad S_{db}^{current} = S_{db} / (\beta_2 \beta_2)$$

$w := w - \alpha \frac{V_{dw}^{current}}{\sqrt{S_{dw}^{current} + \epsilon}}, \quad b := b - \alpha \frac{V_{db}^{current}}{\sqrt{S_{db}^{current} + \epsilon}}$

Combine Momentum and RMSProp

α : needs to be tuned

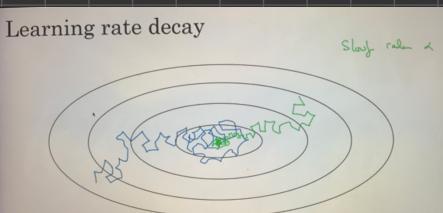
$\beta_1: 0.9$ (δw)

$\beta_2: 0.999$ (δw^2)

$\epsilon: 10^{-8}$

Adam: Adaptive Moment Estimation

Learning rate decay



1 epoch \rightarrow 1 pass through data

$$\alpha_t = \frac{1}{1 + \text{decay_rate} \cdot \text{epoch} + \text{min}}$$

Epoch	α_t
1	0.1
2	≈ 0.09
3	≈ 0.08
4	≈ 0.07
⋮	⋮

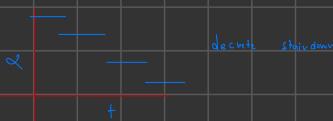
$$\alpha_0 = 0.2$$

decay rate $\rightarrow 1$

Others

$\alpha > 0.95$ equilibrium $\alpha_0 \rightarrow$ Exponentially decay

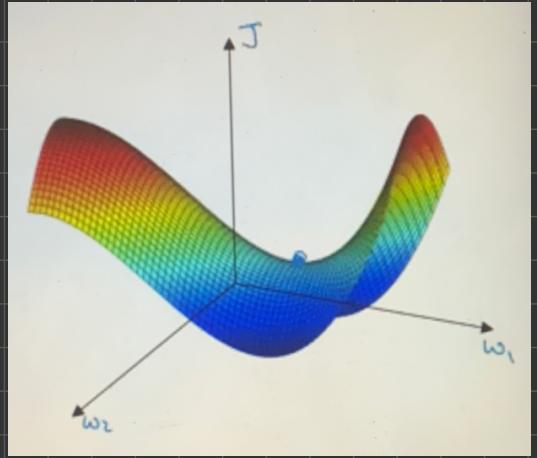
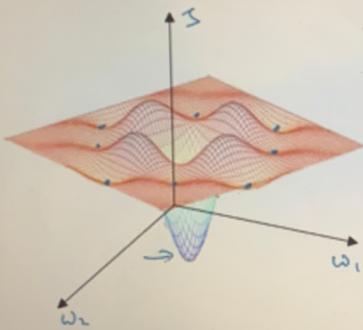
$\alpha = \frac{K}{\text{epoch num}} \cdot \alpha_0 \text{ or } \frac{K}{\sqrt{t}} \cdot \alpha_0$



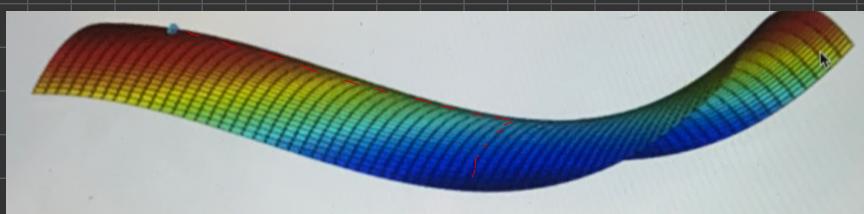
Manual decay

The problem of local optima

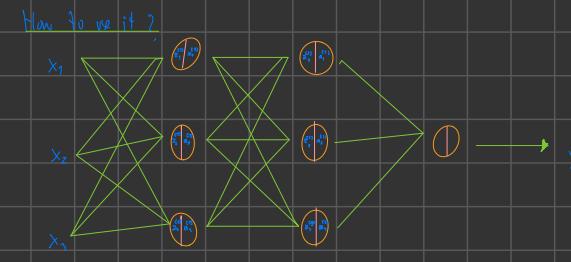
Local optima in neural networks



Platesaus ← The real problems



- Unlikely to get stuck in bad local optima
- Plateaus can make leaving slow



$$x^{t-1}, b^{t-1} \xrightarrow{\text{BN}} z^{t-1} \xrightarrow{f^{t-1}, g^{t-1}} z^{t-1} \rightarrow a^{t-1}, g^{t-1}(z^{t-1}) \xrightarrow{\text{BN}} z^t \xrightarrow{f^t, g^t} z^t \rightarrow a^t, \dots$$

Parameter: $W^{11}, b^{11}, W^{12}, b^{12}, \dots, W^{L1}, b^{L1}$
 $\beta^{t-1}, \gamma^{t-1}, \theta^{t-1}, \phi^{t-1}, \dots, \beta^L, \gamma^L$
 not same as Momentum or RMSprop or Adam

Working with mini-batch

$$X^{t-1} \xrightarrow{W^{t-1}, b^{t-1}} z^{t-1} \xrightarrow{\beta^{t-1}, \gamma^{t-1}} z^{t-1} \rightarrow g^{t-1}(z^{t-1}) \cdot a^{t-1} \xrightarrow{W^t, b^t} z^t \rightarrow \dots$$

$$X^{t-1} \xrightarrow{W^{t-1}, b^{t-1}} z^{t-1} \xrightarrow{\beta^{t-1}, \gamma^{t-1}} z^{t-1} \rightarrow g^{t-1}(z^{t-1}) \cdot a^{t-1} \xrightarrow{W^t, b^t} z^t \rightarrow \dots$$

When we batch training can omit that

$$\begin{aligned} \text{Parameter: } & W^{t-1}, \beta^{t-1}, \gamma^{t-1} \\ & z^{t-1}, W^{t-1} a^{t-1} - b^{t-1} \\ & z^{t-1}, W^{t-1} a^{t-1} \\ & z^{t-1}_{\text{sum}} \\ & z^{t-1}_{\text{sum}} \cdot f^{t-1} z^{t-1}_{\text{sum}} - \mu^{t-1} \end{aligned}$$

with GD

for t-1: Mini-batch epoch

Compute forward pass on X^{t-1}

In each hidden layer use BN to replace z^t s with \hat{z}^{t-1} s

use backprop to compute $dW^{t-1}, d\beta^{t-1}, d\gamma^{t-1}$

Update parameter $W^{t-1} := W^{t-1} - \alpha dW^{t-1}$

$$\beta^{t-1} := \beta^{t-1} - \alpha d\beta^{t-1}$$

$$\gamma^{t-1} := \dots$$

Can use with Adam, momentum, RMSprop

Why does it work

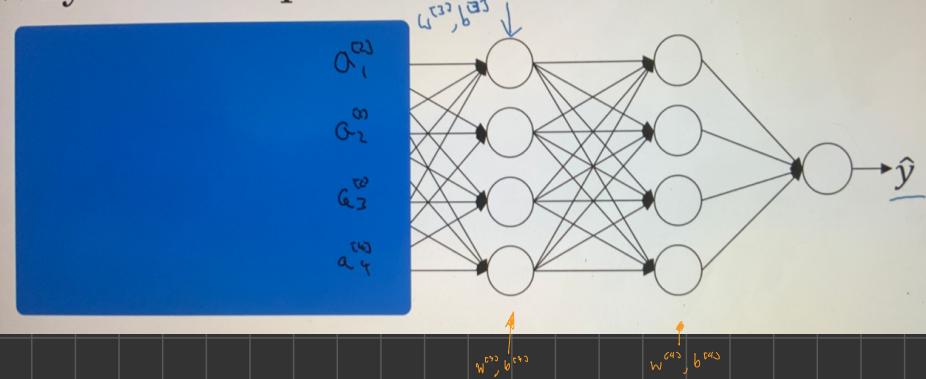


Cat	Non-Cat
$y = 1$	$y = 0$

"Coordinate shift"
 $x \rightarrow \hat{x}$

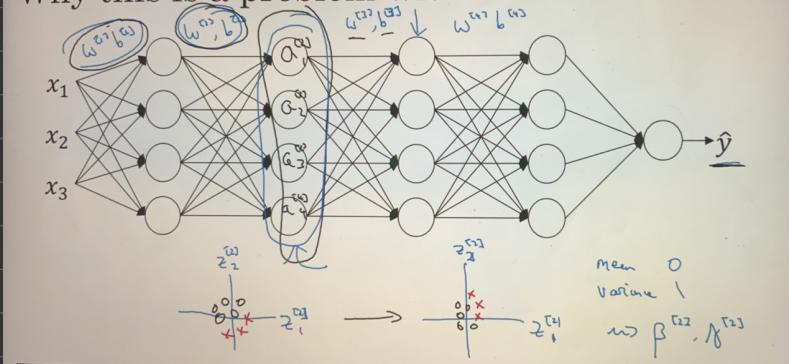
if you've learned some X to Y mapping,

Why this is a problem with neural networks?



Batch Normalization try to help deep layers get to clean mean and variance even the value of deeper layer is update that help deeper layer has ground

Why this is a problem with neural networks?



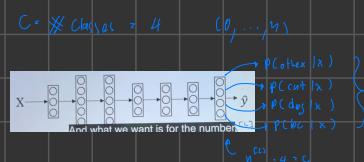
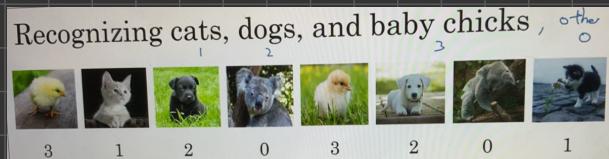
Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

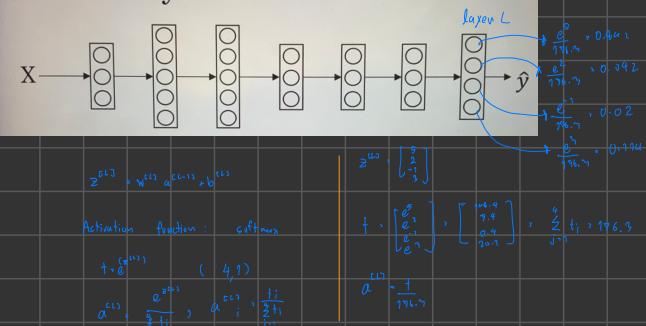
$$\begin{aligned} \mu &= \frac{1}{m} \sum_{i=1}^m x^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2 \\ z_{\text{norm}}^{(i)} &= \frac{x^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ \hat{z}^{(i)} &= \gamma z_{\text{norm}}^{(i)} + \beta \end{aligned}$$

μ, σ^2 : estimate only exponentially with enough (enough) mini-batch

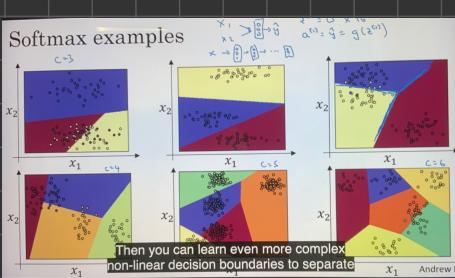
Softmax Regression



Softmax layer



More example



How to use it

$$z^{(i)} = \begin{bmatrix} 1 \\ z_1^{(i)} \\ z_2^{(i)} \\ \vdots \\ z_n^{(i)} \end{bmatrix} \rightarrow \hat{y} = \begin{bmatrix} e^{z_1^{(i)}} \\ e^{z_2^{(i)}} \\ \vdots \\ e^{z_n^{(i)}} \end{bmatrix}$$

"soft max"

"hard max"

$$\hat{y}^{(i)} = g^{(i)}(z^{(i)}) = \frac{e^{z_1^{(i)}}}{\sum_{j=1}^n e^{z_j^{(i)}}} = \begin{bmatrix} 0.842 \\ 0.047 \\ 0.002 \\ \vdots \\ 0.111 \end{bmatrix}$$

Softmax regression generalizes logistic regression to C classes

If $C=2$, softmax reduces to logistic regres.

Loss function

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} - \text{cat}$$

$$a^{(i)}, \hat{y} = \begin{bmatrix} 0.8 \\ 0.1 \\ 0.1 \\ \vdots \\ 0.0 \end{bmatrix}$$

$$\underbrace{\delta(y, \hat{y})}_{\text{loss}} = \sum_j y_j \ln \hat{y}_j$$

$-y \ln \hat{y}_j + -\ln \hat{y}_j$ make $\hat{y} \approx 1$ as possible

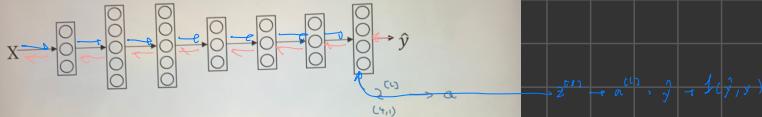
$$J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m \delta(\hat{y}^{(i)}, y^{(i)})$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \quad \hat{Y} = [\hat{y}^{(1)}, \dots, \hat{y}^{(m)}]$$

$$\begin{bmatrix} 0 & 0 & 0 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \quad \begin{bmatrix} 0.8 & \dots \\ 0.1 & \dots \\ 0.1 & \dots \\ \vdots & \ddots \end{bmatrix}$$

$$(m, n) \quad (n, m)$$

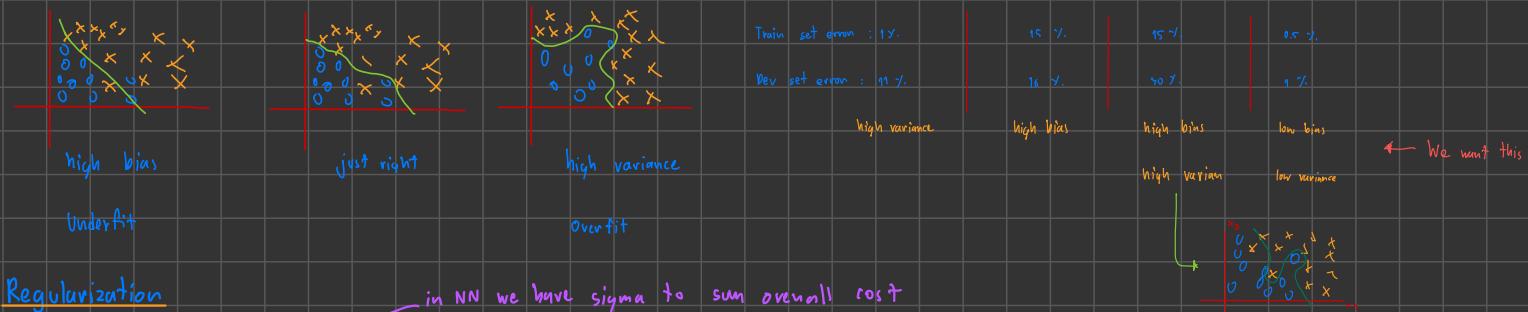
Gradient descent with softmax



$$\text{Backprop} : \frac{d \hat{y}^{(2)}_i}{w_{ij}} = \frac{\partial J}{\partial \hat{y}^{(2)}_i}$$

Summary Optimized NN

Bias / Variance



Regularization

in NN we have sigma to sum overall cost

$$J(w, b) = \sum_{m=1}^M \sum_{i=1}^n \delta(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

omit b λ : regularization parameter

L_2 regularization $\|w\|_F^2 = \sum_{j=1}^n w_j^2 \rightarrow w^T w \xrightarrow{\text{NN}} \|W^{(L)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n (w_{ij}^{(L)})^2$ $w: (n^{(L)}, n^{(L-1)})$ "Frobenius norm"

L_1 regularization $\frac{\lambda}{m} \sum_{l=1}^L \frac{1}{2n} \|w^{(l)}\|_1$

$d_w: \text{Back prop } + \frac{\lambda}{m} w^{(L)}$, $w^{(L)} := w^{(L)} - \alpha d_w^{(L)}$

(In regularization we call weight decay)

Dropout regularization

random prob of layer in NN to delete node to help NN learn faster (but maybe can worse if set too high)

Example

$L=3$ keep_prob = 0.6

$ds = np.random.rand(ds.shape[0], ds.shape[0]) < keep_prob$

$as = np.multiply(as, ds)$

$as /= \text{keep_prob}$ → to not change expect value

Normalizing

Weight Initialization

- Max-min Norm $w^{(L)} = np.random.uniform(shape) * np.sqrt(\frac{1}{n^{(L-1)}})$
- Z-Norm "Xavier initiaton"

Gradient checking

- help to make sure our backprop is correct

Mini batch

Momentum

RMS Prop

Adam

help NN to learn faster

$$\begin{aligned} V_{dw} &= \beta_1 V_{dw} + (1-\beta_1) dw \\ V_{db} &= \beta_1 V_{db} + (1-\beta_1) db \\ S_{dw} &= \beta_2 S_{dw} + (1-\beta_2) dw^2 \\ S_{db} &= \beta_2 S_{db} + (1-\beta_2) db^2 \\ V_{dw}^{corr} &= V_{dw} / (1-\beta_1^t), V_{db}^{corr} = V_{db} / (1-\beta_1^t) \\ S_{dw}^{corr} &= S_{dw} / (1-\beta_2^t), S_{db}^{corr} = S_{db} / (1-\beta_2^t) \\ w &:= w - \frac{\alpha}{\sqrt{S_{dw}^{corr}}} V_{dw}^{corr}, b := b - \frac{\alpha}{\sqrt{S_{db}^{corr}}} V_{db}^{corr} \end{aligned}$$

$\alpha: \text{need to tune}$

Learning rate decay

$\alpha = 0.001 \cdot \text{epoch}^{-0.25}$

$\alpha = \frac{1}{\sqrt{\text{epoch}}}$

Batch Norm

Normalization in hidden

layer not only input

(it help deep layer to

have same mean and variance

even it value is change, like to have ground on it)

Softmax Regression

it change activite

function of output $a^{(i)} = \frac{e^{z^{(i)}}}{\sum_i e^{z^{(i)}}}$

to possibility

$$\begin{bmatrix} 0.9 \\ 0.1 \\ 0.5 \end{bmatrix} \xrightarrow{\text{SM}} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$\beta_1: 0.9$ (dw) Momentum

$\beta_2: 0.999$ (dw²) RMSProp

$\epsilon: 10^{-8}$