

## DL2: Linear regression model

Linear model:  $Y = X\beta + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$

$Y$  = dependent var / outcome var / response var

$X$  =  $p$ -dim independent var, predictor var, explanatory var

$\beta$  =  $p$ -dim regression coefficients:  $\beta_i$ : the effect on the dependent var when increasing the  $i$ th independent var by 1 unit, holding all other  $(p-1)$  predictors constant.

$\varepsilon$  = random variable (error) assumed to be normal.

Assume restricted (linear) form

$$Y = f(X_1, \dots, X_p) + \varepsilon \rightarrow Y = \underbrace{X_1\beta_1 + X_2\beta_2 + \dots + X_p\beta_p}_{\text{fixed part}} + \underbrace{\varepsilon}_{\text{random part}}$$

$E[Y] = X\beta$ ,  
but for fixed  $X^*$ ,  $Y$  differs from  $E[Y]$  by a random  $\varepsilon$

Parameter estimation  $\rightarrow$  obtained by the principle of least squares.

minimize error function:  $\left\{ \begin{matrix} Y \\ n \times 1 \end{matrix} = \begin{matrix} X \\ n \times p \end{matrix} \begin{matrix} \beta \\ p \times 1 \end{matrix} + \begin{matrix} \varepsilon \\ n \times 1 \end{matrix} \right\}$

$$J(\beta) = (y - X\beta)'(y - X\beta) = y'y - 2\beta'X'y + \beta'X'X\beta$$

$$\text{normal } \nabla J(\beta) = -2X'y + 2X'X\beta = 0$$

$$\hat{\beta} = (X'X)^{-1}X'y$$

$$\begin{aligned} \text{cov}(\hat{\beta}) &= \text{cov}((X'X)^{-1}X'y) = (X'X)^{-1}X' \text{cov}(y) = (X'X)^{-1}X'\sigma^2(X'X)^{-1} = \sigma^2(X'X)^{-1} \\ &\approx \hat{\sigma}^2(X'X)^{-1} \end{aligned}$$

standard error of  $\hat{\beta}_i$ :

$$\sqrt{[\hat{\sigma}^2(X'X)^{-1}]_{i,i}} \leftarrow \text{diagonal elements}$$

Error sum of squares (SSE):

$$SSE = (y - \hat{y})'(y - \hat{y}) = \hat{\varepsilon}'\hat{\varepsilon}_{\text{residual}}$$

Coefficient of determination

$$r^2 = 1 - \frac{SSE}{SST}, \quad SST = (y - \bar{y})(y - \bar{y})'$$

Estimated  $\sigma^2$  (variance of error):

$$\hat{\sigma}^2 = \frac{(y - \hat{y})'(y - \hat{y})}{n - p - 1} = \frac{SSE}{n - p - 1}$$

Categorical Independent variable:

$\rightarrow$  set up  $k-1$  dummy variables for  $k$  levels

Hypothesis testing: Omnibus test

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$$

$$H_1: \text{at least one } \beta_k \neq 0$$

Source	df	Sum of Squares	Mean Squares
Regression	$p$	$SSR = (\hat{y} - \bar{y})'(\hat{y} - \bar{y})$	$SSR/p$
Error	$n - p - 1$	$SSE = (y - \hat{y})'(y - \hat{y})$	$SSE/(n - p - 1)$
Total	$n - 1$	$SST = (y - \bar{y})'(y - \bar{y})$	

Table: ANOVA table

$$F = \frac{SSR/p}{SSE/(n - p - 1)} = \frac{MSR}{MSE} \sim F_{\alpha, p, n-p-1}$$

$\uparrow$   
significance level

F-test for subsets of independent variables

$$\text{Full: } y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \varepsilon$$

$$\text{Reduced: } y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

Hypothesis:

$$\begin{aligned} H_0: & \text{reduced model with } \overset{2}{I} \text{ predictors} \\ H_1: & \text{full model with } \overset{4}{p} \text{ predictors} \end{aligned} \quad \left| \begin{array}{c} \text{ANOVA} \\ \rightarrow \end{array} \right. \quad F = \frac{(SSE_R - SSE_F)/(p - I)}{SSE_F/(n - p - 1)} \sim F_{d, p-I, n-p-1}$$

Hypothesis tests for individual regression coefficients

$$H_0: \hat{\beta}_i = 0$$

$$H_1: \hat{\beta}_i \neq 0$$

$$T = \frac{\hat{\beta}_i}{\text{se}(\hat{\beta}_i)} \sim T_{d, p-I, n-p-1}$$

standard error of  $\hat{\beta}_i$ :

$$\sqrt{[\hat{\sigma}^2(X'X)^{-1}]_{i,i}} \leftarrow \sqrt{\text{var}(\hat{\beta}_i)}$$

## DL4: Regularization: Ridge regression & Lasso

assessing whether  $\hat{f}(x) = x'\hat{\beta}$  is a good model:

1)  $\hat{\beta}$  close to  $\beta$ ?  $\Rightarrow \hat{\beta} - \beta = \text{bias}$

2) Will  $\hat{f}(x)$  fit new observations well? (prediction)

### 1) Bias

• MSE for  $\hat{\beta}$

$$\text{MSE}(\hat{\beta}) = E[\|\hat{\beta} - \beta\|^2] = E[(\hat{\beta} - \beta)'(\hat{\beta} - \beta)]$$

### 2) Prediction

• good fitting for current data doesn't guarantee good prediction for new data

small prediction errors (MSPE)

$$\begin{aligned}\text{MSPE}(x_0) &= E[(\hat{f}(x_0) - Y)'(\hat{f}(x_0) - Y)] \\ &= \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0)) + \sigma^2, \quad \text{Bias}(\hat{f}(x_0)) = E[\hat{f}(x_0)] - E[Y]\end{aligned}$$

↳ bias-variance trade off.

- complex model (more terms included)

• Bias  $\downarrow$  (Bias( $\hat{f}(x_0)$ )  $\downarrow$ )

• variance  $\uparrow$  (more sensitive to new data)

$\Rightarrow$  Introducing small bias can lead to substantial decrease in variance (MSPE  $\downarrow$ )

### Example

case ①:  $y \leftarrow (x_{11}, \dots, x_{1000})$

Bias  $\downarrow \Rightarrow$  more independent var  $\rightarrow \downarrow \text{SSE}$

Variance  $\uparrow \Rightarrow$  adding too many var  $\rightarrow$  variance  $\uparrow$  for new, unobserved point.  
 $\rightarrow$  less sensitive to new data

case ②:  $y \leftarrow (x_1, x_2, x_3)$

Bias  $\uparrow$

variance  $\downarrow \Rightarrow$  more robust as only using 3 variables.

By Ridge or Lasso  $\Rightarrow$  Bias  $\uparrow$ , variance  $\downarrow$

$$\Delta \text{Bias} < \Delta \text{variance} \Rightarrow \downarrow \text{MSPE}$$

## 1. Ridge regression & regularization

• unconstrained  $\beta$  can lead to high variance  $\Rightarrow$  regularize  $\beta$  to control variance

• obtain ridge coefficients with constraints:

$$\text{objective functions: } \min(y - X\beta)'(y - X\beta) \text{ s.t. } \sum_{j=1}^p \beta_j^2 \leq t$$

• equivalent to minimizing  $L_2$  loss function:

$$J_2(\beta) = (y - X\beta)'(y - X\beta) + \lambda \|\beta\|_2^2$$

$$\nabla J_2(\beta) = -2X'y + 2XX'\beta + 2\lambda\beta = 0$$

$$\hat{\beta}^{\text{ridge}} = (X'X + \lambda I_p)^{-1} X'y$$

# Neural Network (NN)

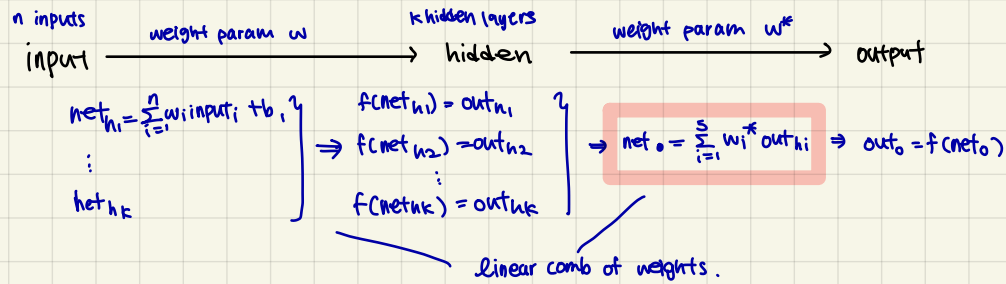
## 1) Forward propagation:

- start with initial weights, calculate output
- Error (actual output - predicted output) calculated

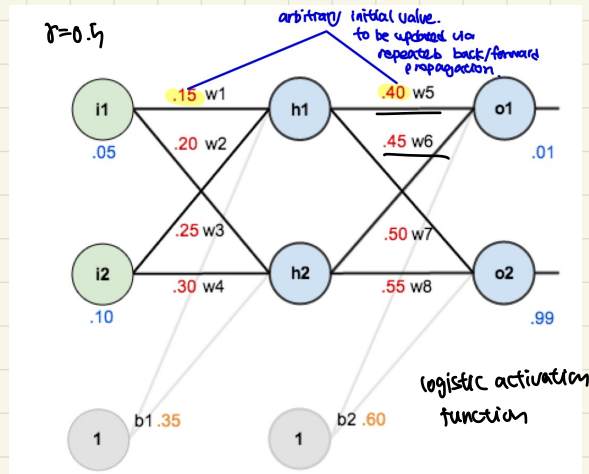
## 2) Backward " :

- weights are updated by reducing the error term in each iteration.
- gradient descent is used

1)



Total error:  $E = \frac{1}{2} (\text{target} - \text{out}_o)^2$



## ① calculation of input $\rightarrow$ hidden layer

$w_1=0.15$   $w_2=0.20$   $w_3=0.25$   $w_4=0.30$   
 $i_1=0.05$   $i_2=0.10$   $b_1=0.35$   $b_2=0.60$

$\text{net}_{h1} = w_{11}i_1 + w_{21}i_2 + b_1 = 0.3775$

$\text{net}_{h2} = w_{31}i_1 + w_{41}i_2 + b_1 = 0.3925$

$\text{out}_{h1} = \frac{1}{1+e^{-\text{net}_{h1}}} = 0.59327$

$\text{out}_{h2} = \frac{1}{1+e^{-\text{net}_{h2}}} = 0.59688$

## ② calculation of hidden layer $\rightarrow$ output

$\text{net}_{o1} = \text{out}_{h1} \cdot w_5 + \text{out}_{h2} \cdot w_6 + b_2 = 1.1059$

$\text{net}_{o2} = \text{out}_{h1} \cdot w_7 + \text{out}_{h2} \cdot w_8 + b_2 = 1.2249$

$\text{out}_{o1} = \frac{1}{1+e^{-\text{net}_{o1}}} = 0.75136$

$\text{out}_{o2} = \frac{1}{1+e^{-\text{net}_{o2}}} = 0.77292$

## ③ calculation of total error

$\Rightarrow \text{out}_{o1} = 0.75136, \text{out}_{o2} = 0.77292$

$O_1 = 0.01, O_2 = 0.99$

$E = \frac{1}{2} (O_1 - \text{out}_{o1})^2 + \frac{1}{2} (O_2 - \text{out}_{o2})^2$

$= 0.29337$

$\Rightarrow$  backpropagation to update weights to minimize E.

2)

- Investigate how much a change in  $w_i$  can affect the total error E:

$$\frac{\partial E}{\partial w_i} = \frac{\partial E}{\partial \text{out}_o} \frac{\partial \text{out}_o}{\partial \text{net}_o} \frac{\partial \text{net}_o}{\partial w_i}$$

- Update weights in the network:

$$w_i^+ = w_i - \gamma \frac{\partial E}{\partial w_i}$$

## ① Backward Pass: output layer. ( $w_5, w_6, w_7, w_8$ )

ex)  $w_5$

$$\frac{\partial E}{\partial w_5} = \underbrace{\frac{\partial E}{\partial \text{out}_{o1}}}_{(1)} \underbrace{\frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}}}_{(2)} \underbrace{\frac{\partial \text{net}_{o1}}{\partial w_5}}_{(3)}$$

(1)  $\frac{\partial E}{\partial \text{out}_{o1}}$   $E = E_{o1} + E_{o2} = \frac{1}{2} (\text{target}_{o1} - \text{out}_{o1})^2 + \frac{1}{2} (\text{target}_{o2} - \text{out}_{o2})^2$

$\frac{\partial E}{\partial \text{out}_{o1}} = -(\text{target}_{o1} - \text{out}_{o1}) = -(0.01 - 0.75136) = +0.74136$

(2)  $\frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}}$   $\text{out}_{o1} = \frac{1}{1+e^{-\text{net}_{o1}}}$

$\frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}} = \frac{e^{-\text{net}_{o1}}}{(1+e^{-\text{net}_{o1}})^2} = \frac{1+e^{-\text{net}_{o1}} - 1}{(1+e^{-\text{net}_{o1}})^2} = \frac{1+e^{-\text{net}_{o1}}}{(1+e^{-\text{net}_{o1}})^2} - \frac{1}{(1+e^{-\text{net}_{o1}})^2}$

$= \frac{1}{1+e^{-\text{net}_{o1}}} - \frac{1}{(1+e^{-\text{net}_{o1}})^2} = \text{out}_{o1} - (\text{out}_{o1})^2 = \text{out}_{o1}(1 - \text{out}_{o1})$

$= 0.75136(1 - 0.75136) = 0.1868$

(3)  $\frac{\partial \text{net}_{o1}}{\partial w_5}$   $\text{net}_{o1} = \text{out}_{h1} \cdot w_5 + \text{out}_{h2} \cdot w_6 + b_2$

$\frac{\partial \text{net}_{o1}}{\partial w_5} = \text{out}_{h1} = 0.59327$

$\therefore \frac{\partial E}{\partial w_5} = 0.74136 \times 0.1868 \times 0.59327 = 0.08216704$

## ② Backward pass: update $w_7$

$w_7^+ = w_7 - \gamma \frac{\partial E}{\partial w_7} = 0.4 - 0.5 \times 0.082167 = 0.3589$

repeat for  $w_6^+, w_7^+, w_8^+$

$\left[ \begin{array}{l} w_5 \\ w_6 \end{array} \right] \rightarrow \text{net}_{o1} \rightarrow \text{out}_{o1} \rightarrow E$   
 $\left[ \begin{array}{l} w_7 \\ w_8 \end{array} \right] \rightarrow \text{net}_{o2} \rightarrow \text{out}_{o2} \rightarrow E$

③ Backward Pass. Hidden layer ( $w_1, w_2, w_3, w_4$ )

ex)  $w_1 \rightarrow \text{net}_{h1} \rightarrow \text{out}_{h1} \rightarrow \text{net}_{o1} \rightarrow \text{out}_{o1} \rightarrow E_{o1} \rightarrow E$   
 $\rightarrow \text{net}_{o2} \rightarrow \text{out}_{o2} \rightarrow E_{o2} \rightarrow E$  } more complicated.

$$\frac{\partial E}{\partial w_1} = \underbrace{\frac{\partial E}{\partial \text{out}_{h1}}}_{(1)} \cdot \underbrace{\frac{\partial \text{out}_{h1}}{\partial \text{net}_{h1}}}_{(2)} \cdot \underbrace{\frac{\partial \text{net}_{h1}}{\partial w_1}}_{(3)}$$

$$(1) \frac{\partial E}{\partial \text{out}_{h1}} = \frac{\partial E_{o1}}{\partial \text{out}_{h1}} + \frac{\partial E_{o2}}{\partial \text{out}_{h1}}$$

$$\begin{aligned} \rightarrow \frac{\partial E_{o1}}{\partial \text{out}_{h1}} &= \frac{\partial E_{o1}}{\partial \text{out}_{o1}} \cdot \frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}} \cdot \underbrace{\frac{\partial \text{net}_{o1}}{\partial \text{out}_{h1}}}_{\text{net}_{o1} = \text{out}_{h1} \cdot w_3 + \text{out}_{h2} \cdot w_4 + b_2} \\ &= 0.74176 \times 0.1868 \times 0.4 \\ &= 0.05539 \end{aligned}$$

$\frac{\partial \text{net}_{o1}}{\partial \text{out}_{h1}} = w_3$

$$\frac{\partial E_{o2}}{\partial \text{out}_{h1}} = \frac{\partial E_{o2}}{\partial \text{out}_{o2}} \cdot \frac{\partial \text{out}_{o2}}{\partial \text{net}_{o2}} \cdot \frac{\partial \text{net}_{o2}}{\partial \text{out}_{h1}}$$

=

$$(2) \frac{\partial \text{out}_{h1}}{\partial \text{net}_{h1}} \quad \text{out}_{h1} = \frac{1}{1 + e^{-\text{net}_{h1}}}$$

$$\frac{\partial \text{out}_{h1}}{\partial \text{net}_{h1}} = \text{out}_{h1} (1 - \text{out}_{h1}) = 0.59327 (1 - 0.59327) = 0.2413$$

$$(3) \frac{\partial \text{net}_{h1}}{\partial w_1} \quad \text{net}_{h1} = w_1 i_1 + w_2 i_2 + b_1$$

$$\frac{\partial \text{net}_{h1}}{\partial w_1} = i_1 = 0.05$$

## Multilayer Neural Networks.

↳ work with complex non-linear representations of data

→ adding more layers & neurons → ↑ specialization to train data, ↓ performance on test data

- perceptron
  - simple perceptron: outputs are binary →  $y = \sum w_i x_i$ . If  $y > 0$ , output = 1  
If  $y < 0$ , output = -1
  - multi-class perceptron: many possible labels.

- linear separable classifier

- if  $f(b + \sum w_i x_i) > \text{threshold}$  → output = 1 (active)
- if  $f(b + \sum w_i x_i) < \text{threshold}$  → output = -1 (inactive).
- If classification is incorrect, change  $w$ .  
else, don't change  $w$ .

- weight updates

- 1) get  $x$  and output label  $y$
- 2) Initialize  $w, b$  for  $f(x)$
- 3) If  $f(x) = y$ , mark as completed; else, fix it
- 4) adjust score based on error.
  - $f(x) = \text{sign}(b + \sum w_i x_i)$
  - if  $y = +1$  and  $f(x) = -1$ ,  $\sum w_i x_i$  is too small, make it bigger  
 $y = -1$  and  $f(x) = +1$ ,  $\sum w_i x_i$  is too large, make it smaller.
- 5)  $w = w$  if  $f(x) = y$ , else  $w = w + yx$
- 6)  $b = b$  if  $f(x) = y$ , else  $b = b + yR$ ,  $R = \max(\|x\|)$
- 7) repeat 3~6 until  $f(x) = y$

## CNN

dim of input datasets  
Input (image) :  $28 \times 28 \times 3$  array of numbers,  $n=1000$  → # of data samples  
resolution RGB values

output: class [cat, dog, ...] → probability of classes [0.29, 0.49, 0.1, 0.09, 0.17]

Structure

Input → Conv → ReLU → Conv → ReLU → Pool → ReLU → Conv → ReLU → Pool → FC → output  
Filter the image extract important info from images through dim reduc. Shrink image size voting & classification. generates vector of prob

1) Convolutional layer.

$32 \times 32 \times 3 \rightarrow 28 \times 28 \times 1$  conv  
one  $5 \times 5 \times 3$  filter, stride=1  
 $\rightarrow 28 \times 28 \times 2$  conv  
two  $5 \times 5 \times 3$ , stride=1  
⇒ multiplications between the filter & pixel values.

↑ stride ⇒ overlaps less ⇒ more dimension reduction

⇒ as we keep applying conv layers, the dim of the volume decreases.

⇒ we want to preserve original input volume in the early layers of NN so that we can extract some info.

⇒ padding to prevent fast volume reduction

\* padding

three  
 $32 \times 32 \times 3 \rightarrow 7 \times 7 \times 3$  filter, but want to keep output volume as  $32 \times 32 \times 3$

⇒ apply zero padding size 2. →  $36 \times 36 \times 3$  input volume.

padding size filter size

$$P = \frac{K-1}{2} \quad (\text{to keep input \& output dim the same})$$

input height/length

$$O = \frac{(W - K + 2P)}{S} + 1$$

stride size

2) ReLU activation function

$$y = xI(x > 0)$$

### 3) Pooling

ex) maxpooling of 2x2 filter with stride 2

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

take max → 6 8 → reduce dim  
14 16

### 4) FC Layer.

→ at the end of CNN

→ takes an input volume → outputs an N-dim vector of prob.

- optimization of CNN

- initial values : randomize filter & filter values (weights)
- tuning param: stride, padding, filter size
- minimize error :  $E = \frac{1}{2} (\text{target} - \text{out})^2$

### Back propagation: CNN

Assume we have

$$\frac{\partial E}{\partial \text{net}} = \frac{\partial E}{\partial \text{out}} \frac{\partial \text{out}}{\partial \text{net}}$$

Backpropagation

•  $\frac{\partial E}{\partial b}$  : how the change in b can affect E

$$b^* = b - \gamma \frac{\partial E}{\partial b}$$

•  $\frac{\partial E}{\partial w}$  : how the change in w " " E

$$w^* = w - \gamma \frac{\partial E}{\partial w}$$

•  $\frac{\partial E}{\partial x}$  : how the change in pixel  $x_{ij}$  can affect E

### 1-Dim Input Example

Input :  $x = (x_1, x_2, x_3, x_4)$

$w = (w_1, w_2, b)$

Net :  $\text{net} = (\text{net}_1, \text{net}_2, \text{net}_3)$

if stride = 1

$$(x_1 \ x_2 \ x_3 \ x_4)$$

$w_1 \ w_2$

$$\begin{aligned} x_1 w_1 + x_2 w_2 + b &= \text{net}_1 \\ x_2 w_1 + x_3 w_2 + b &= \text{net}_2 \\ x_3 w_1 + x_4 w_2 + b &= \text{net}_3 \end{aligned}$$

applying 1 kernel,  
output volume of  
x1 dim

$$\frac{\partial E}{\partial b} = \frac{\partial E}{\partial \text{net}} \frac{\partial \text{net}}{\partial b}$$

$$= \left[ \frac{\partial E}{\partial \text{net}_1} \quad \frac{\partial E}{\partial \text{net}_2} \quad \frac{\partial E}{\partial \text{net}_3} \right] \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \frac{\partial E}{\partial \text{net}_1} + \frac{\partial E}{\partial \text{net}_2} + \frac{\partial E}{\partial \text{net}_3}$$

$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial \text{net}} \frac{\partial \text{net}}{\partial w} = \begin{bmatrix} x_1 & x_2 \\ x_2 & x_3 \\ x_3 & x_4 \end{bmatrix}$$

$$= \left[ \frac{\partial E}{\partial \text{net}_1} \quad \frac{\partial E}{\partial \text{net}_2} \quad \frac{\partial E}{\partial \text{net}_3} \right] \begin{bmatrix} x_1 & x_2 \\ x_2 & x_3 \\ x_3 & x_4 \end{bmatrix}$$

$$= \left[ \frac{\partial E}{\partial \text{net}_1} x_1 + \frac{\partial E}{\partial \text{net}_2} x_2 + \frac{\partial E}{\partial \text{net}_3} x_3 \quad \frac{\partial E}{\partial \text{net}_2} x_2 + \frac{\partial E}{\partial \text{net}_3} x_3 + \frac{\partial E}{\partial \text{net}_3} x_4 \right] = \begin{bmatrix} \frac{\partial E}{\partial w_1} \\ \frac{\partial E}{\partial w_2} \end{bmatrix}$$

$$\frac{\partial E}{\partial x} = \left[ \frac{\partial E}{\partial \text{net}_1} \quad \frac{\partial E}{\partial \text{net}_2} \quad \frac{\partial E}{\partial \text{net}_3} \right] \begin{bmatrix} w_1 & w_2 & 0 & 0 \\ 0 & w_1 & w_2 & 0 \\ 0 & 0 & w_1 & w_2 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} \frac{\partial E}{\partial x_1} \\ \frac{\partial E}{\partial x_2} \\ \frac{\partial E}{\partial x_3} \\ \frac{\partial E}{\partial x_4} \end{bmatrix} = \begin{bmatrix} \frac{\partial E}{\partial \text{net}_1} w_1 \\ \frac{\partial E}{\partial \text{net}_1} w_2 + \frac{\partial E}{\partial \text{net}_2} w_1 \\ \frac{\partial E}{\partial \text{net}_2} w_2 + \frac{\partial E}{\partial \text{net}_3} w_1 \\ \frac{\partial E}{\partial \text{net}_3} w_2 \end{bmatrix}$$

2dim input example.

Input:

$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & & & \\ x_{31} & \dots & & \\ x_{41} & & & \end{bmatrix}$$

$$W = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$$

$b$

Output

$$net = \begin{bmatrix} net_{11} & net_{12} & net_{13} \\ net_{21} & net_{22} & net_{23} \\ net_{31} & net_{32} & net_{33} \end{bmatrix}$$

<conv with one filter  $w$ , stride=1, padding=0>

$$net_{11} = w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22} + b$$

$$net_{12} = w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23} + b$$

$\vdots$

$\Rightarrow 3 \times 3$  output volume

$$* net_{ij} = \left( \sum_{k=1}^2 \sum_{l=1}^2 w_{kl} x_{i+k-1, j+l-1} \right) + b, \quad \forall i, j \in \{1, 2, 3\}$$

## RNN

↳ unlike feed forward networks, RNNs can use internal memory for their processing.

↳ RNN takes the previous output or hidden states as inputs.

The input at time  $t$  has some historical information at time  $T < t$ .

↳ Intermediate values (state) of RNNs can save info about past inputs

### Equations

- Hidden nodes:  $a_t = W_H h_{t-1} + W_X x_t$    
 num of hidden nodes  $d \times 1$   $d \times d$   $d \times 1$   $d \times d$   $d \times 1$    
  $2d \times d^2 \rightarrow$  total # of weight param that you have to estimate.
- Output from hidden state:  $h_t = \tanh(a_t)$    
  $d \times 1$   $d \times 1$   $d \times d$   $d \times 1$
- Prediction at time  $t$ :  $y_t = \text{softmax}(W_y h_t)$    
  $d \times 1$   $d \times d$   $d \times 1$

↳ we only train RNN one time from each randomized weights  $W_X, W_Y, W_H$ .

Repeat optimization via backpropagation.  $\Rightarrow W_X, W_Y, W_H$  will change.

RNN updates

$$W_X^+ = W_X - \gamma \frac{\partial L}{\partial W_X}$$

$$W_Y^+ = W_Y - \gamma \frac{\partial L}{\partial W_Y}$$

$$W_H^+ = W_H - \gamma \frac{\partial L}{\partial W_H}$$

Multi-class Cross entropy Loss function.

$$L(y, \hat{y}) = -y' \log(\hat{y})$$

Binary ex)

$$y = (0.99, 0.01) \quad \hat{y}_1 = (0.99, 0.01) \quad \hat{y}_2 = (0.01, 0.99)$$

$$\text{loss for } \hat{y}_1: L(y, \hat{y}_1) = -y' \log(\hat{y}_1)$$

$$L_{\text{total}}(y, \hat{y}) = L_1(y_1, \hat{y}_1) + L_2(y_2, \hat{y}_2) + L_3(y_3, \hat{y}_3)$$

chainrule for updating  $W_Y$ .

$$z_t = W_Y h_t \rightarrow \hat{y}_t = \text{softmax}(z_t) \rightarrow L_t = -y_t' \log(\hat{y}_t)$$

$$\text{chain rule: } \frac{\partial L_t}{\partial W_Y} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial W_Y}$$

$$\text{gradient: } \frac{\partial L_{\text{total}}}{\partial W_Y} = \sum_t \frac{\partial L_t}{\partial W_Y}$$

$$\text{update: } W_Y^+ = W_Y - \gamma \frac{\partial L_{\text{total}}}{\partial W_Y}$$

### Chain rule for updating $W_X$

$$h_t = \tanh(W_H h_{t-1} + W_X x_t) \rightarrow z_t = W_Y h_t \rightarrow \hat{y}_t = \text{softmax}(z_t) \rightarrow L_t = -y_t' \log(\hat{y}_t)$$

chainrule:

$$\frac{\partial L_t}{\partial W_X} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_X} + \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_X} + \dots$$

$$\frac{\partial L_t}{\partial W_X} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_X} = \sum_{k=0}^t \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_X} \cdot \frac{\partial h_k}{\partial W_X}$$

$$\text{gradient: } \frac{\partial L_{\text{total}}}{\partial W_X} = \sum_{t=1}^n \sum_{k=0}^t \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_X} \cdot \frac{\partial h_k}{\partial W_X}$$

$$\text{update: } W_X^+ = W_X - \gamma \frac{\partial L_{\text{total}}}{\partial W_X}$$

### chain rule for updating $W_H$

$$h_t = \tanh(W_H h_{t-1} + W_X x_t) \rightarrow z_t = W_Y h_t \rightarrow \hat{y}_t = \text{softmax}(z_t) \rightarrow L_t = -y_t' \log(\hat{y}_t)$$

chain rule:

$$\frac{\partial L_t}{\partial W_H} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_H} + \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_H} \cdot \frac{\partial h_{t-1}}{\partial W_X} + \dots$$

### RNN issue: vanishing gradient (+ deep NNS)

: weight update requires chain rule which has lots of components.

so if each component close to 0, gradient becomes 0  $\Rightarrow$  no update.

↳ happens because we assume all the previous informations  $(h_1, \dots, h_n)$  gives impact on  $h_t$ .

↳ solution: LSTM.

LSTM: have gates that make the long/short term memory to be stored.

forget gate & input gate decides which information to be abandoned & stored.