

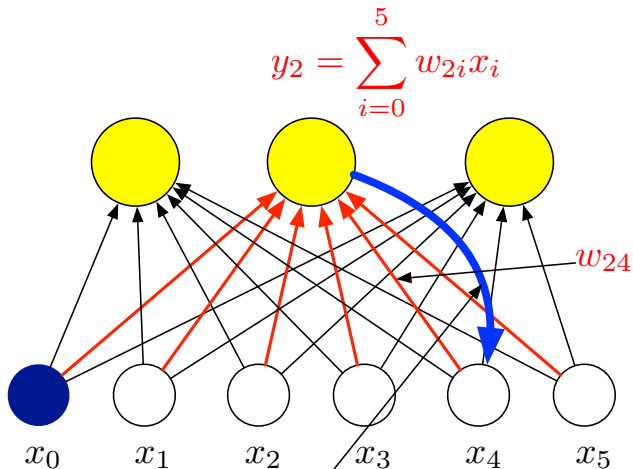
Stochastic gradient descent; Classification

Steve Renals

Machine Learning Practical — MLP Lecture 2
30 September 2015

Single Layer Networks

Applying gradient descent to a single-layer network



$$\Delta w_{24} = \sum_n (y_2^n - t_2^n) x_4^n$$

Example: Rainfall Prediction

Daily Southern Scotland precipitation (mm). Values may change after QC.

Alexander & Jones (2001, Atmospheric Science Letters).

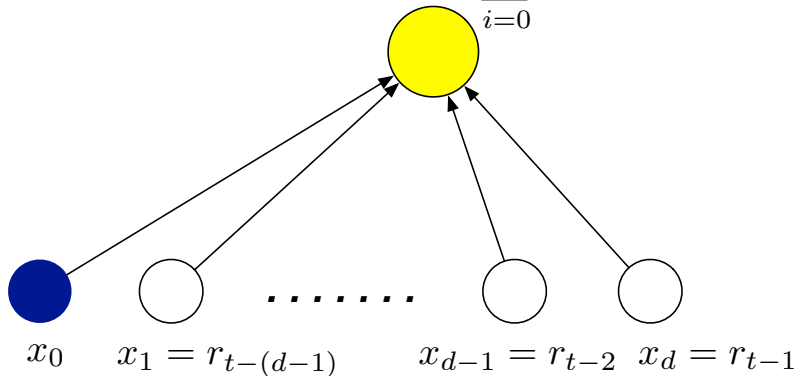
Format=Year, Month, 1-31 daily precipitation values.

1931	1	1.40	2.10	2.50	0.10	0.00	0.00	0.90	6.20	1.90	4.90	7.30	0.80	0.30	2
1931	2	0.90	0.60	0.40	1.10	6.69	3.00	7.59	7.79	7.99	9.59	24.17	1.90	0.20	4
1931	3	0.00	1.30	0.00	0.00	0.00	0.50	0.40	0.60	1.00	0.00	0.10	7.30	6.20	0
1931	4	3.99	3.49	0.00	2.70	0.00	0.00	1.80	1.80	0.00	0.20	3.39	2.40	1.40	1
1931	5	1.70	0.00	0.70	0.00	5.62	0.70	13.14	0.80	11.13	11.23	0.60	1.70	10.83	8
1931	6	1.40	16.40	3.70	0.10	5.80	12.90	4.30	4.50	10.40	13.20	0.30	0.10	9.30	29
1931	7	9.49	1.70	8.69	4.10	2.50	13.29	2.70	5.60	3.10	1.30	7.59	3.90	2.30	7
1931	8	0.20	0.00	0.00	0.00	0.00	0.60	2.00	0.60	6.60	0.60	0.90	1.20	0.50	4
1931	9	9.86	4.33	1.01	0.10	0.30	1.01	0.80	1.31	0.00	0.30	4.23	0.00	1.01	1
1931	10	23.18	5.30	4.20	6.89	4.10	11.29	10.09	5.80	11.99	1.80	2.00	5.10	0.30	0
1931	11	6.60	20.40	24.80	3.30	3.30	2.60	5.20	4.20	8.00	13.60	3.50	0.90	8.50	15
1931	12	3.20	21.60	16.00	5.80	8.40	0.70	6.90	4.80	2.80	1.10	1.10	0.90	2.50	3
1932	1	12.71	41.12	22.51	7.20	12.41	5.70	1.70	1.80	24.41	3.80	0.80	13.71	4.30	17
1932	2	0.00	0.22	0.00	0.54	0.33	0.11	0.00	0.00	0.22	0.11	0.22	0.00	0.00	0
1932	3	0.10	0.00	0.00	1.60	8.30	4.10	10.00	1.10	0.00	0.00	0.00	0.60	0.50	0
1932	4	7.41	4.61	1.10	0.10	9.41	8.61	2.10	13.62	17.63	4.71	0.70	0.30	10.02	3
1932	5	0.10	0.20	0.00	0.10	0.70	0.10	0.80	1.00	0.30	0.00	10.51	17.42	4.11	1
1932	6	0.00	0.00	0.00	0.20	0.00	0.00	0.60	0.20	0.50	0.00	0.00	0.10	0.00	0
1932	7	2.41	7.62	13.94	7.42	1.30	1.30	1.80	3.81	2.61	4.01	1.00	4.81	9.93	0
1932	8	0.00	1.70	0.30	1.00	2.70	4.61	3.40	2.60	0.50	1.30	9.61	1.80	3.81	0
1932	9	19.37	7.39	9.69	2.70	3.50	3.79	16.68	5.29	4.69	16.88	3.50	1.00	14.08	2
1932	10	4.40	0.50	0.10	1.80	6.40	8.20	14.69	18.39	4.30	2.80	0.10	16.19	2.20	0
1932	11	11.37	8.08	5.79	0.00	0.00	0.00	0.00	0.20	0.00	0.00	0.10	0.30	0.00	0
1932	12	20.23	19.93	3.81	2.40	0.00	0.00	0.00	0.10	0.40	0.40	0.10	0.70	2.30	13
1933	1	3.40	28.50	2.80	18.80	5.30	4.50	14.60	8.80	0.60	3.50	0.00	3.10	0.50	19
1933	2	6.10	2.60	14.80	33.10	8.00	9.00	3.10	4.70	7.00	0.10	0.10	0.90	0.10	0
1933	3	2.59	5.29	3.99	5.99	7.19	7.09	0.30	29.54	5.19	0.00	0.00	0.00	1.10	3
1933	4	0.40	14.98	3.20	0.50	0.00	0.00	0.00	11.98	1.70	0.10	4.69	0.20	0.00	0
1933	5	0.00	0.00	4.71	9.92	2.21	13.73	3.81	5.71	1.80	0.10	0.80	0.20	0.00	0

Single Layer Network for Rainfall Prediction

Output - predicted observation

$$y = \hat{r}_t = \sum_{i=0}^d w_i x_i$$



Input - previous $d-1$ observations

Stochastic Gradient Descent (SGD)

- Training by batch gradient descent is *very slow* for large training data sets
 - The algorithm sums the gradients over the entire training set before making an update
 - Since the update steps (η) are small many updates are needed
- Solution: **Stochastic Gradient Descent (SGD)**
- In SGD the true gradient $\partial E / \partial w_{ki}$ (obtained by summing over the entire training dataset) is approximated by the gradient for a point $\partial E^n / \partial w_{ki}$
- The weights are updated after each training example rather than after the batch of training examples
- Inaccuracies in the gradient estimates are washed away by the many approximations
- To prevent multiple similar data points (all with similar gradient approximation inaccuracies), present the training set in random order

SGD Pseudocode (linear network)

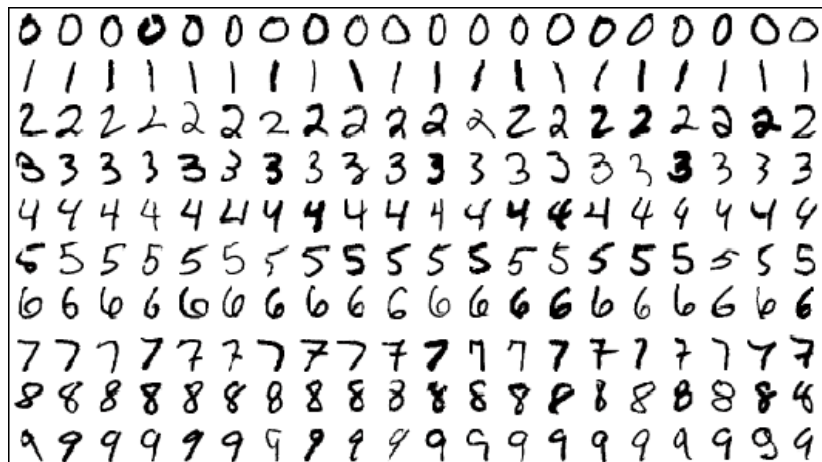
```
1: procedure SGDTRAINING(X, T, W)
2:   initialize W to small random numbers
3:   randomize order of training examples in X
4:   while not converged do
5:     for  $n \leftarrow 1, N$  do
6:       for  $k \leftarrow 1, K$  do
7:          $y_k^n \leftarrow \sum_{i=0}^d w_{ki} x_i^n$ 
8:          $\delta_k^n \leftarrow y_k^n - t_k^n$ 
9:         for  $i \leftarrow 1, d$  do
10:           $w_{ki} \leftarrow w_{ki} - \eta \cdot \delta_k^n \cdot x_i^n$ 
11:        end for
12:      end for
13:    end for
14:  end while
15: end procedure
```

Minibatches

- Batch gradient descent – compute the gradient from the batch of N training examples
- Stochastic gradient descent – compute the gradient from 1 training example each time
- Intermediate – compute the gradient from a **minibatch** of M training examples – $M > 1$, $M \ll N$
- Benefits of minibatch:
 - Computationally efficient by making best use of vectorisation, keeping processor pipelines full
 - Possibly smoother convergence as the gradient estimates are less noisy than using a single example each time

Classification

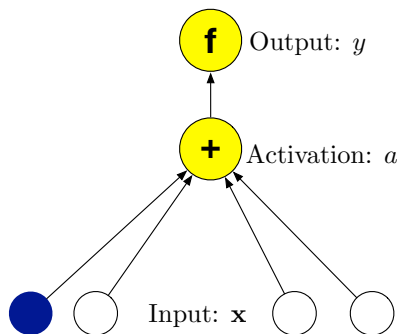
MNIST Digit Classification



Classification and Regression

- **Regression:** predict the value of the output given an example input vector - e.g. what will be tomorrow's rainfall (in mm)
- **Classification:** predict the category given an example input vector - e.g. will it be rainy tomorrow (yes or no)?
- Classification outputs:
 - **Binary:** 1 (yes) or 0 (no)
 - **Probabilistic:** p , $1 - p$ (for a 2-class problem)
- One could train a linear single layer network as a classifier:
 - Output targets are 1/0 (yes/no)
 - At run time if the output $y > 0.5$ classify as yes, otherwise classify as no
- This will work, but we can do better....
- Output activation functions to constrain the outputs to binary or probabilistic (logistic / sigmoid)

Two-class classification

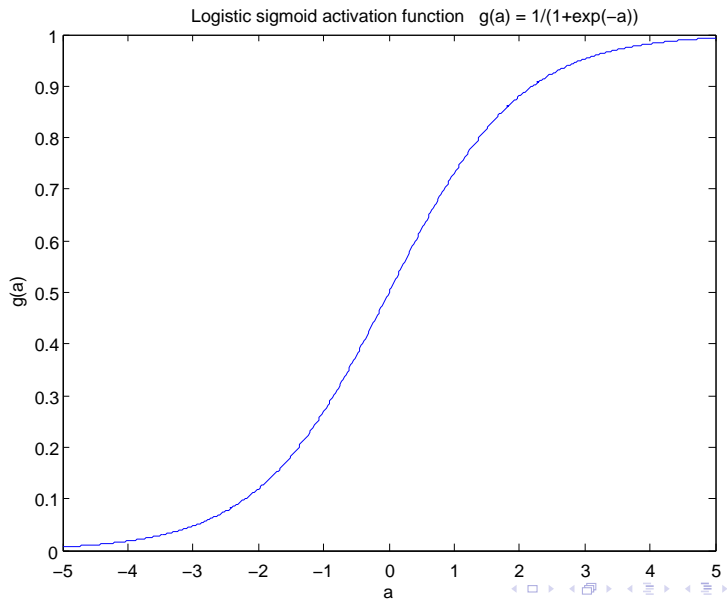


Single-layer network, binary/sigmoid output

$$\text{Binary (step function): } f(a) = \begin{cases} 1 & \text{if } a \geq 0.5 \\ 0 & \text{if } a < 0.5 \end{cases}$$

$$\text{Probabilistic (sigmoid function): } f(a) = \frac{1}{1 + \exp(-a)}$$

Sigmoid function



Sigmoid single layer networks

- Binary output: activation is not differentiable. Can use *perceptron learning* to train binary output single layer networks
- Probabilistic output: sigmoid single layer network (statisticians would call this logistic regression). Let a be the *activation* of the single output unit, the value of the weighted sum of inputs, before the activation function, so:

$$y = f(a) = f\left(\sum_i w_i x_i\right)$$

- Two classes, so single output y , with weights w_i

Sigmoid single layer networks

- Training sigmoid single layer network: Gradient descent requires $\partial E/\partial w_i$ for all weights:

$$\frac{\partial E^n}{\partial w_i} = \frac{\partial E^n}{\partial y^n} \frac{\partial y^n}{\partial a^n} \frac{\partial a^n}{\partial w_i}$$

For a sigmoid:

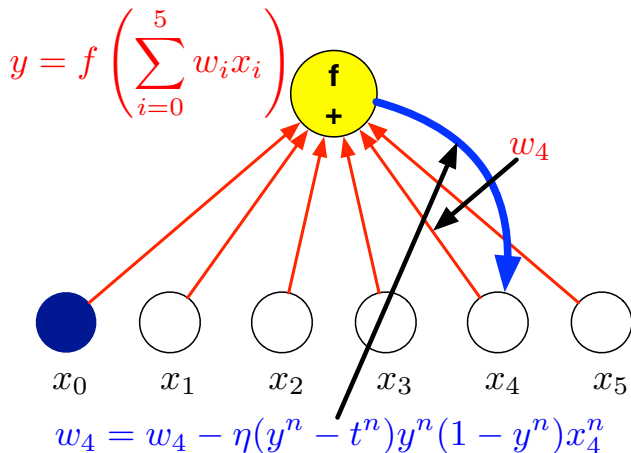
$$y = f(a) \quad \frac{dy}{da} = f(a)(1 - f(a))$$

(Show that this is indeed the derivative of a sigmoid.)

- Therefore:

$$\frac{\partial E^n}{\partial w_i} = \underbrace{(y^n - t^n)}_{\delta^n} \underbrace{f(a^n)(1 - f(a^n))}_{f'(a^n)} x_i^n$$

Applying gradient descent to a sigmoid single-layer network



Cross-entropy error function (1)

- If we use a sigmoid single layer network for a two class problem (C_1 (target $t = 1$) and C_2 ($t = 0$)), then we can interpret the output as follows

$$y \sim P(C_1 | \mathbf{x}) = P(t = 1 | \mathbf{x})$$
$$(1 - y) \sim P(C_2 | \mathbf{x}) = P(t = 0 | \mathbf{x})$$

- Combining, and recalling the target is binary

$$P(t | x, \mathbf{W}) = y^t \cdot (1 - y)^{1-t}$$

This is a Bernoulli distribution. We can write the log probability:

$$\ln P(t | x, \mathbf{W}) = t \ln y + (1 - t) \ln(1 - y)$$

Cross-entropy error function (2)

- Optimise the weights \mathbf{W} to maximise the log probability – or to minimise the negative log probability. Write the error function as follows:

$$E^n = -(t^n \ln y^n + (1 - t^n) \ln(1 - y^n)) .$$

This is called the **cross-entropy error function**

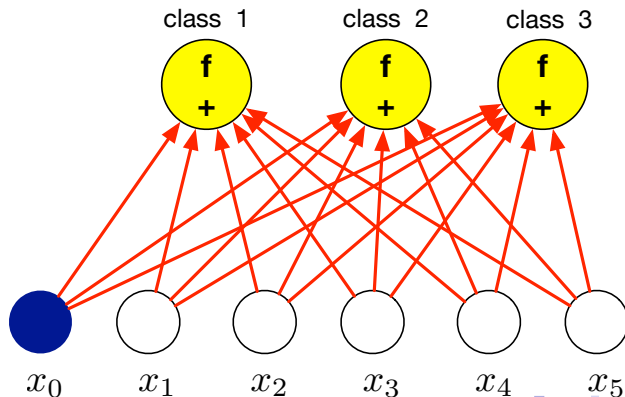
- Gradient descent training requires the derivative $\partial E / \partial w_i$ (where w_i connects the i th input to the single output).

$$\begin{aligned} \frac{\partial E}{\partial y} &= -\frac{t}{y} + \frac{1-t}{1-y} = \frac{-(1-y)t + y(1-t)}{y(1-y)} \\ \frac{\partial E}{\partial w_i} &= \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial a} \cdot \frac{\partial a}{\partial w_i} \\ &= \frac{-(1-y)t + y(1-t)}{y(1-y)} \cdot y(1-y) \cdot x_i = (y-t)x_i \end{aligned}$$

(derivative of the sigmoid $y(1-y)$ cancels)

Multi-class networks

- If we have K classes use a “one-hot” (“one-from-N”) output coding – target of the correct class is 1, all other targets are zero
- It is possible to have a multi-class net with sigmoids



Multi-class networks

- If we have K classes use a “one-hot” (“one-from-N”) output coding – target of the correct class is 1, all other targets are zero
- It is possible to have a multi-class net with sigmoids
- This will work... but we can do better
- Using multiple sigmoids for multiple classes means that $\sum_k P(k|\mathbf{x})$ is not constrained to equal 1 – we want this if we would like to interpret the outputs of the net as class probabilities
- Solution – an activation function with a sum-to-one constraint: **softmax**

$$y_k = \frac{\exp(a_k)}{\sum_{j=1}^K \exp(a_j)}$$
$$a_k = \sum_{i=0}^d w_{ki} x_i$$

- This form of activation has the following properties
 - Each output will be between 0 and 1
 - The denominator ensures that the K outputs will sum to 1
- Using softmax we can interpret the network output y_k^n as an estimate of $P(k|\mathbf{x}^n)$
- Softmax is the multiclass version of the two-class sigmoid

Softmax – Training (1)

- We can extend the cross-entropy error function to the multiclass case

$$E^n = - \sum_{k=1}^C t_k^n \ln y_k^n$$

- Again the overall gradient we need is

$$\begin{aligned} \frac{\partial E^n}{\partial w_{ki}} &= \sum_{c=1}^C \frac{\partial E}{\partial y_c} \cdot \frac{\partial y_c}{\partial a_k} \cdot \frac{\partial a_k}{\partial w_{ki}} \\ &= \sum_{c=1}^C -\frac{t_c}{y_c} \cdot \frac{\partial y_c}{\partial a_k} \cdot x_i \end{aligned}$$

Note that the k th activation a_k – and hence the weight w_{ki} – influences the error function through all the output units, because of the normalising term in the denominator. We have to take this into account when differentiating.

Softmax – Training (2)

- Do the differentiation, and you will find

$$\frac{\partial y_c}{\partial a_k} = y_c(\delta_{ck} - y_k)$$

Where δ_{ck} ($\delta_{ck} = 1$ if $c = k$, $\delta_{ck} = 0$ if $c \neq k$) is called the Kronecker delta

- And putting it all together

$$\frac{\partial E^n}{\partial w_{ki}} = (y_k^n - t_k^n)x_i^n$$

The delta rule!

- 1 Modify the SGD pseudocode for sigmoid outputs
- 2 Modify the SGD pseudocode for softmax outputs
- 3 Modify the SGD pseudocode for minibatch
- 4 For softmax and cross-entropy error, show that

$$\frac{\partial E^n}{\partial w_{ki}} = (y_k^n - t_k^n)x_i^n$$

(use the quotient rule of differentiation)

Summary

- Stochastic gradient descent (SGD) and minibatch
- Classification and regression
- Sigmoid activation function and cross-entropy
- Multiple classes – Softmax
- Next lecture: multi-layer networks and hidden units