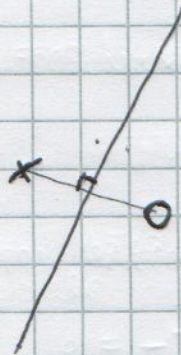
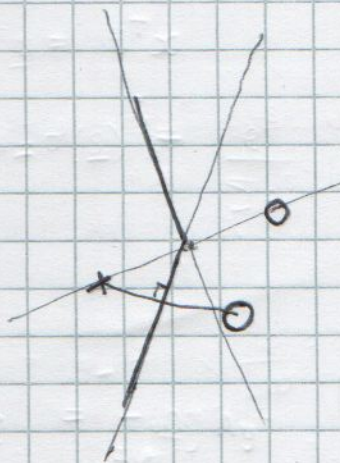


6.1



2 points  
1 ar bisector



add more bisectors.

6.1.9

want

$|x_i - w_j|^2$  is  $\begin{cases} \text{small, } c_i = j \\ \text{large, } c_i \neq j \end{cases}$

$$|x_i - w_j|^2 = |x_i|^2 + |w_j|^2 - 2w_j^T x_i$$

$$\text{suppose } |x_i|^2 = |w_j|^2 = 1$$

$$\text{then } \underline{w_j^T x_i} = \begin{cases} \text{large, } c_i = j \\ \text{small, } c_i \neq j \end{cases}$$

6.2

$$h = w^T x$$

fit by L2  $w^* = \arg \min_w e = |h - y|^2 = |w^T x - y|^2$

$$\frac{\partial e}{\partial w} = 2(w^T x - y) x$$

$$\frac{\partial e}{\partial w} = 0, \quad x(x^T w - y) = 0 \quad x x^T w - x y = 0$$

many  $x$ 's  $\sum_i x_i x_i^T w - \sum_i y_i x_i = 0$

let  $X = \begin{pmatrix} \cdot & x_1^T & \cdot \\ \cdot & x_2^T & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix}$

$$X^T X w - X^T Y = 0$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \end{pmatrix}$$

$$w = (X^T X)^{-1} X^T Y$$

$w$  vector same size as data  $x$

$$h = \frac{w^T x}{\text{3072x1 vectors}} + b_0$$

3072x1  
vectors

bias

$$\underline{h} = w^T X + \underline{b}$$

$$\begin{bmatrix} x \\ 1 \end{bmatrix}$$

6.3

$$h = w^T x$$

$$\begin{pmatrix} s_1 \\ s_2 \\ s_3 \\ \vdots \end{pmatrix} = \begin{pmatrix} w_1^T \\ w_2^T \\ w_3^T \\ \vdots \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \end{pmatrix}$$

$$s_i = w_i^T x$$

$w_i$  are "templates" for class  $i$

6.4

$$y = a_0 + a_1 x + a_2 x^2 + a_3 x^3$$

$y_i, x_i$

4 unknowns, solve for  $a_0, a_1, a_2, a_3$  using 4 points

$N^{\text{th}}$  order poly  $N+1$  points

6.5

$$y = Mx$$

L2 loss  $e = |y - Mx|^2$

regularization  $e = |y - Mx|^2 + \lambda |a|^2$

6.7

$$e = \frac{1}{2} |w^T x - t|^2$$

$$\frac{\partial e}{\partial w} = (w^T x - t) x^T$$

update rule  $w_{t+1} = w_t - \alpha \frac{\partial e}{\partial w} = \nabla w$

$$w_{t+1} = w_t - \alpha (w^T x - t) x^T$$

learning rate

6.8

SGD + momentum

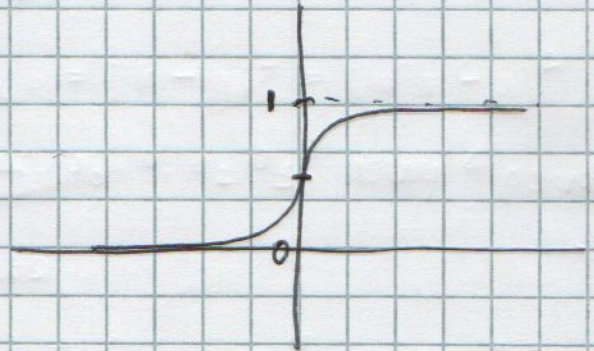
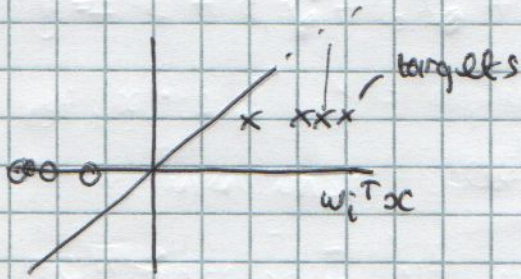
$$v_{t+1} = \rho v_t + \nabla w_t$$

$$\rho = 0.9 - 0.99$$

$$w_{t+1} = w_t - \alpha v_{t+1}$$

improved versions Ada grad, rms prop, Adam.

6.9



$$h(s_i) = \frac{1}{1 + e^{-s}} \quad \text{logistic}$$

$$\underline{h}(s) = \frac{e^{-s}}{\sum_i e^{s_i}}$$

$$\underline{s} = w^T x$$

~~6.9~~

6.10

define scores  $\underline{s} = \underline{w}^T \underline{x}$

softmax predictor  $h = \sigma(\underline{s}) = \frac{e^{s_i}}{\sum_j e^{s_j}}$       $\sigma_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$

cross entropy loss

$$e = - \sum_i t_i \log h_i = - \log h_y$$

$$\frac{\partial e}{\partial \underline{s}} = - \frac{1}{h_y} \frac{\partial}{\partial \underline{s}} h_y = - \frac{1}{h_y} \frac{\partial}{\partial \underline{s}} \frac{e^{s_y}}{\sum_i e^{s_i}}$$

$$\frac{\partial}{\partial s_k} \frac{e^{s_y}}{\sum_i e^{s_i}} = \begin{cases} -\frac{e^{s_y}}{(\sum_i e^{s_i})^2} e^{s_k} & k \neq y \\ \frac{e^{s_y}}{\sum_i e^{s_i}} - \frac{e^{s_y} e^{s_k}}{(\sum_i e^{s_i})^2} & k = y \end{cases}$$

$$= \begin{cases} -h_y h_k, & k \neq y \\ h_y - h_y^2, & k = y \end{cases}$$

$$\frac{\partial e}{\partial \underline{s}} = - \frac{1}{h_y} \frac{\partial}{\partial \underline{s}} h_y = \begin{cases} h_k, & k \neq y \\ -1 + h_k, & k = y \end{cases}$$

$$\frac{\partial e}{\partial \underline{s}} = \underline{h} - \underline{t}$$

where  $\underline{h} = \sigma(\underline{w}^T \underline{x})$       $\underline{s} = \underline{w}^T \underline{x}$

$e$  is cross entropy loss.

$$\frac{\partial e}{\partial \underline{w}} = \frac{\partial e}{\partial \underline{s}} \frac{\partial \underline{s}}{\partial \underline{w}} = \underline{(\underline{h} - \underline{t})} \underline{x}^T$$