Last time: we started with logistic regression.

$$p(y=1|x) = \sigma(w^T x + w_0)$$

Aside: one way to derive this: let's have the difference in $p(y=1|x)$ and $p(y=0|x)$ be linear in log space:

$$\ln p(y=1|x) - \ln p(y=0|x) = \boxed{w^T x + w_0}$$

Now, we substitute $p(y=0|x) = (1 - p(y=1|x))$

$$w^T x + w_0 = \ln \frac{p(y=1|x)}{p(y=0|x)} = \ln \frac{p(y=1|x)}{1 - p(y=1|x)}$$

and solve for $p(y=1|x)$:

$$\frac{p(y=1|x)}{1 - p(y=1|x)} = \exp\{w^T x + w_0\}$$

$$p(y=1|x) = \frac{\exp\{w^T x + w_0\}}{1 + \exp\{w^T x + w_0\}} = \frac{1}{1 + \exp\{-1 \cdot \boxed{(w^T x + w_0)}\}}$$

$$\downarrow$$
$$f(x, w_0, w_1)$$

Write down the log likelihood:

$$ll(w) = \sum_{n=1}^{N} -y_n \ln(1 + \exp\{-f_n\}) - (1-y_n) \ln(1 + \exp\{f_n\}) \qquad \text{to max.}$$

$$\underbrace{\qquad}$$
$$f_n = w^T x_n + w_0$$

(came from
$$p(y=1|x)^y$$
$$y \ln p(y=1|x)$$

$$\mathcal{L}(w) = \sum_{n=1}^{N} y_n \ln(1 + \exp\{f_n\}) + (1-y_n) \ln(1 + \exp\{f_n\}) \qquad \text{to min}$$

$$\searrow$$

Now we can take gradients $\rightarrow$ going to use chain rule

define
$$\mathcal{L}_n(w) =$$
$$y_n \ln(1 + \exp\{-f_n\})$$
$$+ (1-y_n) \ln(1 + \exp\{f_n\})$$

$$\frac{\partial \mathcal{L}(w)}{\partial w} = \frac{\partial f_n}{\partial w} \frac{\partial \mathcal{L}_n(w)}{\partial f_n}$$

$$\Downarrow \qquad \Downarrow$$

$$\frac{\partial}{\partial w} w^T x_n + w_0$$
$$= x_n$$

$$\frac{\partial}{\partial f_n} \ln(1 + \exp\{f_n\}) = \frac{\exp\{f_n\}}{1 + \exp\{f_n\}} = p(y_n=1|x_n)$$

$$\frac{\partial}{\partial f_n} \ln(1 + \exp\{-f_n\}) = \frac{-\exp\{-f_n\}}{1 + \exp\{-f_n\}} = -p(y_n=0|x_n)$$

$$\text{so:} \quad \frac{\partial L_n(\omega)}{\partial f_n} = y_n \, p(y_n=0|x_n) \dots$$

$$\frac{\partial L_n(\omega)}{\partial \omega} = (x_n)\Big( y_n \, p(y_n=0|x_n)(-1) + (1-y_n)\, p(y_n=1|x_n)\Big)$$

so now, we can apply (S)GD:

     1: sample a batch of data (size = M)    $(x_m, y_m)$

     2: compute $\dfrac{\partial L_m(\omega)}{\partial \omega}$ for all $x_m, y_m$ in M

     3: finally, $\quad \omega \leftarrow \omega - \eta \dfrac{N}{M} \sum_m \dfrac{\partial L_m(\omega)}{\partial \omega}$   ] a lot of additional innovations here to reduce variance

                    step size   const           → ADAM (dim-wise step size)

                                                → momentum

     assumption: our errors add. (log likelihoods, $\sum err_n^2$, ..)

are we done?

     well, $\omega^T x + \omega_0$ isn't the most expressive function ...

     choose a basis! $[x_1 \dots x_D] \rightarrow [\phi_1(\underline{x}) \cdots \phi_J(\underline{x})]$

         ↪ great if you know what $\phi$ should be

         ↪ but choosing $\phi$ can be hard

next: can we <u>learn</u> $\phi$, alongside $\omega, \omega_0$?        ASIDE
(Neural Networks are <u>one</u> way to do this!)     "Adaptive Basis" in stats Regression

**◻ Neural Networks**

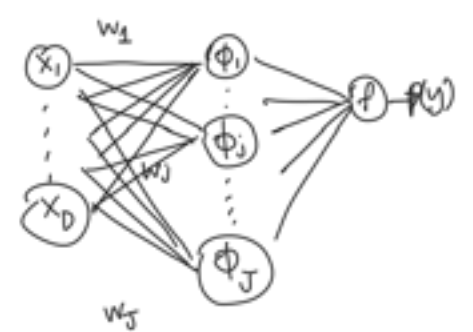     idea: suppose we <u>nest</u> a LR in a LR!

         $f(x_n) = \omega^T \phi + \omega_0$    } linear ; will go through sigmoid

                                $p(y=1|x_\bullet) = \sigma(f_n)$

         $\phi_j = \sigma(w_j^T x + w_{j0})$ for $j$ in $1 \dots J$
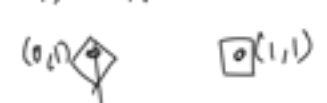
in pictures:



easy to put into matrix form:

if I take all the $w_j$'s: $\quad [\,w_1^1, w_2^1 \dots w_j^1 \dots w_J^1\,] \}D$

                                           $W^1$

$$\phi = \sigma\Big(W^{1^T} X \oplus W_0^1\Big)$$

applied pone-wise    $J \times N$    vector of size $J$

examples (toy) suppose we want an X-OR function:

     (0,0)◇      ◻(1,1)     <u>not</u> linearly separable!

(0,0) (1,0)

We can create basis functions that are of the $\sigma$ form to make this linearly separable!

let's "pick out" [0,0] $\qquad \phi_1 = \sigma(-x_1 - x_2 + \frac{1}{2})$

$\Rightarrow \quad f = -\phi_1 - \phi_2 + \frac{1}{2}$

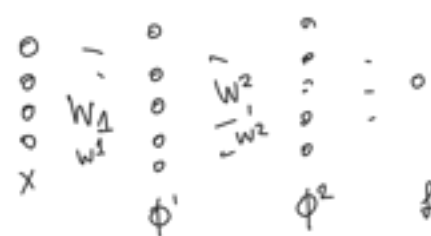let's "pick out" [1,1] $\qquad \phi_2 = \sigma(x_1 + x_2 - 1.5)$

Another example:



how could a linear set of basis functions be a good classifier for this?

pos: class is this intersection in the middle

Coming up:



Deep neural networks

multiple layers; other architectures

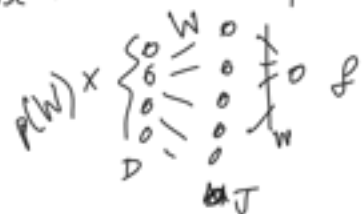multiple outputs for multiple classes via softmax:

○ $P(y=1|x)$
○ :
○ :
○ $P(y=k|x)$

$$\frac{\exp\{w_k^T x\}}{\sum_\ell \exp\{v_\ell^T x\}}$$

Concept Exercise : Bayesian Neural Networks

BNNs: let's put prob. over $W, W^1$ & compute posteriors...

suppose we have a 1-layer network w/ $J$-nodes



1. let's write down $f$ as a function of feature outputs $\phi_j(x)$, $w_j \leftarrow w^T x$

$$f = \sum_j w_j \phi_j(x)$$

let $\phi_j(x)$ have some mean $\mu_j$ and variance $\sigma_j^2$ $\qquad$ var(xg) = var(x)var(y) +

2. let
and $w_m \sim N(0, \sigma_w^2)$

$$\text{var}(x)\mu_y^2 +$$
$$\text{var}(y)\mu_x^2$$

what is the mean & variance of $f$?

$$E[f] = 0 \qquad \sigma_f^2 = J\left[\sigma_w^2 \sigma_J^2 + \sigma_w^2 \mu_J^2\right] \qquad \text{assuming all units same } \mu_J, \sigma_J^2$$

3. Can you argue that $f$ is Gaussian for large $J$?

CLT

↳ can be used to That an only broad 1-layer NN w/ these priors

............ (dist over functions)