$$\underline{L = 3}$$

$$\hat{f}(\underline{x}) = \hat{f}_3(\hat{f}_2(\hat{f}_1(\underline{x})))$$



$$h^{(1)} = g^{(1)}(W^{(1)}\underline{x} + b^{(1)})$$

$$\underset{M_1 \times P}{} \quad \underset{R^P}{} \quad \underset{M_1}{}$$

$$h^{(2)} = g^{(2)}(W^{(2)}h^{(1)} + b^{(2)})$$

$$\underset{M_2 \times M_1}{} \quad \underset{M_1}{} \quad \underset{M_2}{}$$

$$\hat{y} = g^{(3)}(W^{(3)}h^{(2)} + b^{(3)}) \longleftarrow \mathbb{R}$$

$$\underset{1 \times M_2}{} \quad \underset{M_2}{} \quad \underset{1}{}$$

# Why are NNs interesting?

**Claim:** NNs are universal approximators.

True, but also true of many other methods.

  e.g. poly regression, KNN

No, mention of convergence rate.

---

NNs are automatic feature engineering machines.

One way to create complex methods is via feature engineering.

Start w/ $\underset{\sim}{x} \in \mathbb{R}^P$

Could fit

$$\hat{f}(x) = \underset{\sim}{x}^T \beta + \beta_0 \quad \text{(linear model]}$$

but I could also create a feature map

$$\phi : \mathbb{R}^P \to \mathbb{R}^M \quad (\text{typ. } M > P)$$

$$\phi : \mathbb{R} \to \mathbb{R}^{M} \quad (\text{typ. } M > P)$$

and then fit a linear model to this:

$$\hat{f}(\underset{\sim}{x}) = \phi(\underset{\sim}{x})^T \beta + \beta_0.$$

<u>Problem</u>: need to specify $\phi$ manually

<u>Instead</u>! Work w/ parameterized collection
of feat maps $\phi_\alpha$
ad fit

$$\hat{f}(\underset{\sim}{x}) = \phi_\alpha(\underset{\sim}{x})^T \beta + \beta_0$$

Then I need to learn parameters
$$(\beta, \beta_0, \alpha)$$

<u>Called</u>: feature learning. — non-lin.

<u>e.g.</u> $\phi_\alpha(\underset{\sim}{x}) = g\left(W\underset{\sim}{x} + b\right)$

when $\alpha = (W, b)$

hidden

This is equiv. to a single ^hidden layer NN:

$$\hat{f}(\underset{\sim}{x}) = \beta_0 + \beta^T g(W\underset{\sim}{x} + b).$$

Deep NNs w/ lots of layers are just learning really complicated feat maps $\phi_\alpha$

$$\phi_\alpha(x) = \hat{f}_{L-1}(\hat{f}_{L-2}(\cdots \hat{f}_2(\hat{f}_1(x))))$$

$$\alpha = \left(W^{(L-1)}, b^{(L-1)}, W^{(L-2)}, b^{(L-2)}, \ldots\right)$$

---

<u>For regression</u>: predicting $y \in \mathbb{R}$

last layer typ uses $g^{(L)}(x) = x$

so that

$$\hat{f}(\underset{\sim}{x}) = \underbrace{W^{(L)}}_{1 \times M_{L-1}} \underbrace{h^{(L-1)}}_{M_{L-1}} + \underbrace{b^{(L)}}_{1}$$

For classification: predict $y \in \{1, ..., K\}$

Typically use $g^{(L)}(h) = \text{softmax}(h)$

↶ K-vector

$$\text{Softmax}(h)_i = \frac{e^{h_i}}{\sum_{j=1}^{K} e^{h_j}}$$

K-vector

$z = \text{Softmax}(h)$

then ① $z_i \geqslant 0$

② $\sum_{i=1}^{K} z_i = 1$ .

$$O(x) = \text{Softmax}\left(W^{(L)} h^{(L-1)} + b^{(L)}\right)$$

then

$$\hat{f}(x) = \underset{k}{\text{argmax}} \; O(x)_k$$

NNs generalize many methods.

NNs generalize many methods:

<u>E.g.</u> $L = 1$, $g^{(1)}(x) = x$ then

$$\hat{f}(x) = W^{(1)}\underset{\sim}{x} + b^{(1)}$$

this is just regression.

<u>E.g.</u> $L = 1$ and

$$O(\underline{x}) = \text{Softmax}(W^{(1)}\underset{\sim}{x} + b^{(1)})$$

$$\hat{f}(x) = \underset{k}{\text{argmax}} \; O(x)_p$$

then this is just logistic regression.

---

For FNNs there are lots of choices in architecture to be made

    ① how many layers (depth)

    ② how many hidden units at each

② how many hidden units at each layer (width)

③ which <u>activation</u> functions
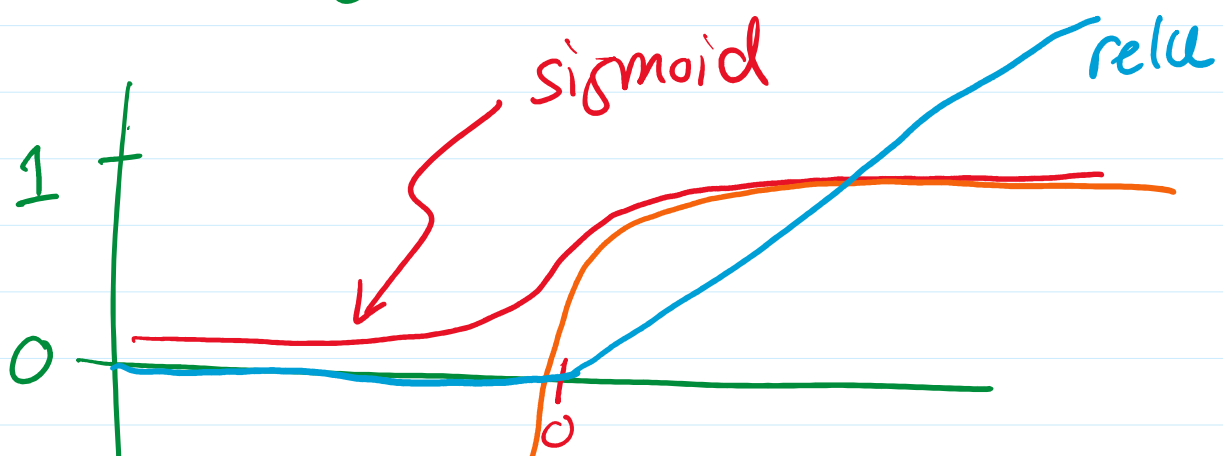
<u>Modern wisdom</u>: deeper networks are better than wider.
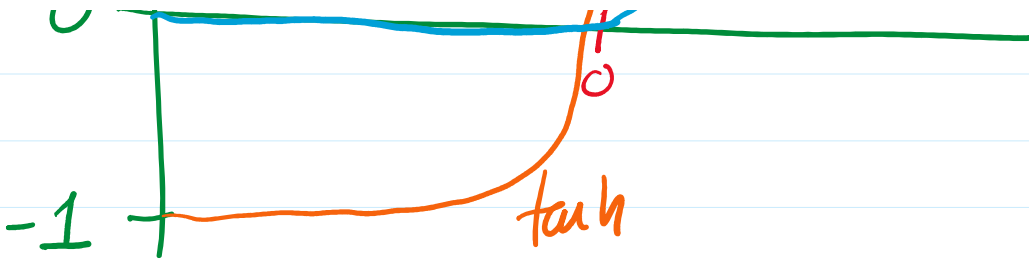
<u>Activation Functions</u> :

  <u>Sigmoid</u> : $g(h) = \dfrac{1}{1 + e^{-h}}$

  <u>tanh</u> : $g(h) = \tanh(h)$

  <u>ReLU</u> : $g(h) = \max(0, h)$

$-1$    tanh

---

## How do we learn the params?

Our params:

$$\Theta = (W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \ldots, W^{(L)}, b^{(L)})$$

this is super high dim'l.

To learn need to specify some loss

$$L(y, f_\theta(\underset{\sim}{x}))$$

and then find the value of $\Theta$ that minimizes the empirical loss

$$\boxed{\hat{\Theta} = \arg\min_{\Theta} \frac{1}{N} \sum_{n=1}^{N} L(y_n, f_\theta(\underset{\sim}{x}_n))}$$

training data of

For regression: $L(y, \hat{y}) = (y - \hat{y})^2$

Squared error

K-class classification:

$$O(\underline{x}) = \text{Softmax}\left(W^{(L)} b^{(L-1)} + b^{(L)}\right)$$

K probs that sum to $1$
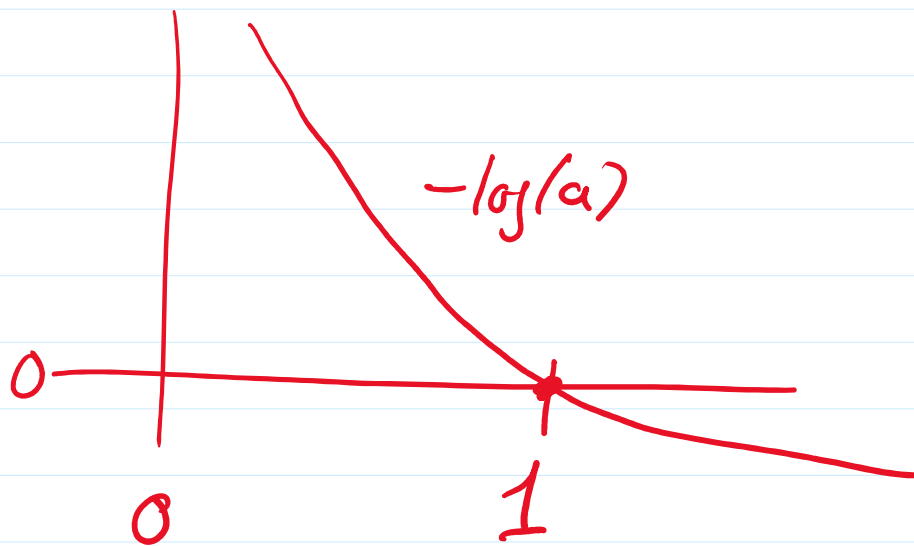(for each class)

$$\tilde{y}_k = \mathbb{I}(y = k)$$

$$\left[ y = 2 \quad : \quad \tilde{y} = (0, 1, 0, 0 \dots ) \right]$$

Use Cross-entropy loss:

$$L(y, O(\underline{x})) = -\sum_{k=1}^{K} \tilde{y}_k \log(O(\underline{x})_k)$$

$$= -\log(O(\underline{x})_y)$$

$-\log(a)$

0

0          1

---

This is difficut b/c

① $\theta$ is super high dim'l

② may hone lots of data

③ $\hat{f}$ is very complex.

---

To optimize use gradient descent

$$\sum^{N} / (\quad \hat{f}(x_i))$$

$$\hat{\Theta} = \underset{\Theta}{argmin} \sum_{n=1}^{N} L(y_n, \hat{f_\Theta}(x_n))$$

$$\mathcal{L}(\theta)$$

<u>Idea</u>: follow the slope of $\mathcal{L}$:

- gradient $\frac{\partial \mathcal{L}}{\partial \Theta}\big|_{\Theta_0}$ points in dir of fastest ascent at $\Theta_0$

- $-\frac{\partial \mathcal{L}}{\partial \Theta}\big|_{\Theta_0}$ points (locally) in dir to move to decrease $\mathcal{L}$ fastest

<u>Grad Desc</u>:

Intialize $\Theta^{(0)}$

For $t = 1, 2, 3, \ldots$

$$\Theta^{(t)} = \Theta^{(t-1)} - \alpha \nabla_\Theta \mathcal{L}\big|_{\Theta^{(t-1)}}$$

$$\theta \sim |\theta^{(t-1)}$$

grad.

step size
(learning rate)

**Problem:** $N$ may be really large and so Calc. grad may be slow

$$\mathcal{L}(\theta) = \sum_n L(y_n, \hat{f_\theta}(x_n))$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = \sum_n \frac{\partial L_n}{\partial \theta}$$

**Soln:** <u>Stochastic Grad Descent (SGD)</u>

Idea: instead of Calc grad over all trains, we just use a (random) subset.

$\longrightarrow$ mini-batch