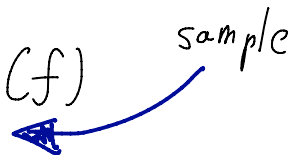# Boosting
# Decision Trees
# Neural Networks

Notes are available
online (see Piazza)

# Boosting, Decision Trees, Neural Networks

* Revisit surrogate loss : exp-loss
* Boosting and exp-loss connected
* Decision trees, exp-loss, and boosting connected
* From boosted classifiers to Neural Networks

Notes:

* old-new view that is not covered in one book or a few papers

* Intuitive, simple, "proof-free" (almost)

* Focus on empirical loss $\mathcal{L}_S(f)$    sample

# Exponential Loss

$f : \mathcal{X} \to \mathbb{R}$     $|f(x)|$ — confidence in prediction

$\text{sign}(f(x))$ — predicted outcome
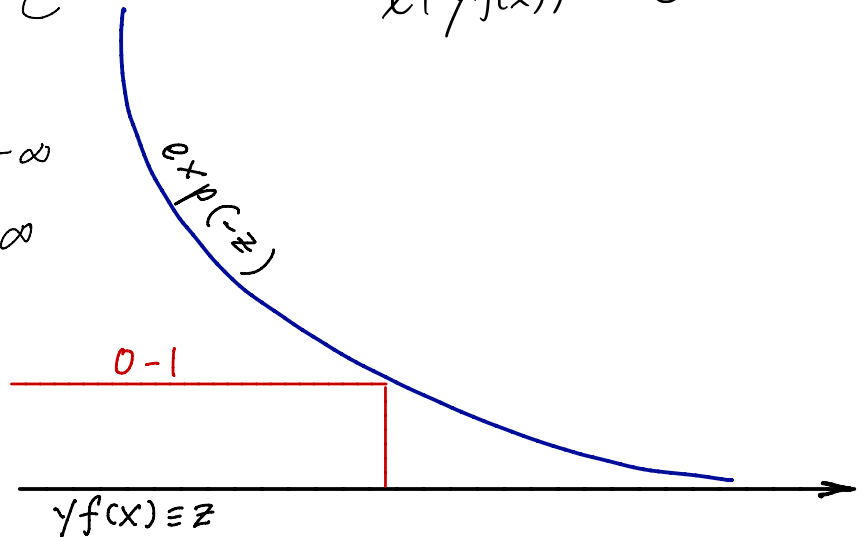
$$\mathcal{L}_{\mathcal{S}}^{exp}(f) \equiv \frac{1}{|S|} \sum_{i=1}^{|S|} e^{-y_i \, f(x_i)} \qquad \ell(y f(x)) = e^{-y f(x)}$$

* grows fast ↗ as $z \to -\infty$
* gets small ↘ as $z \to \infty$
* "nice" properties

exp(-z)

0-1

$y f(x) \equiv z$

# Predictors which can abstain

$f(x) \in \{-1, 0, +1\}$

$f: \mathcal{X} \longrightarrow \{-1, 0, +1\}$

no
negative

IDK
IDC

yes
positive

## Example

$x \in \mathbb{R}^d \quad w \in \mathbb{R}^d$

$$f(x) = \begin{cases} -1 & w \cdot x < -\sigma \\ 0 & |w \cdot x| \leq \sigma \\ +1 & w \cdot x > \sigma \end{cases}$$

Yes

$w \cdot x = 0$

??

N

$w$

$\sigma$

$\sigma$

# Calibrating Predictors

* Assume that $f: x \to \{-1, 0, +1\}$ is given us

* Need to find $\alpha$ such that $\mathcal{L}_g^{exp}(\alpha f(x))$

is minimized

Can use $\alpha_+, \alpha_-$

* Process called calibration / rescaling / reweighing

$$\min_{\alpha \in \mathbb{R}} \frac{1}{|S|} \sum_{i=1}^{|S|} e^{-y_i (\alpha f(x_i))}$$

$$\tilde{x}_i \in \mathbb{R}$$

$$\underline{\text{Search for } \alpha}$$

* Exploit the fact that $\quad f: \mathcal{X} \to \{-1, 0, +1\}$

$$S_+ = \{i \mid y_i \, f(x_i) = +1\} \quad \text{Correct}$$

$$S \quad \longrightarrow \quad S_0 = \{i \mid y_i \, f(x_i) = 0\} \quad \text{Blank}$$

$$\boxed{|S| = m} \quad \longrightarrow \quad S_- = \{i \mid y_i \, f(x_i) = -1\} \quad \text{Mistake}$$

$$\mathcal{L}_S^{exp}(\alpha f) = \frac{1}{m}\left( \sum_{i \in S_+} e^{-\alpha} + \sum_{i \in S_-} e^{+\alpha} + \sum_{i \in S_0} e^{0} \right)$$

Are we in a better shape?

# Closed form solution

Define $\quad \mu_+ = \dfrac{|S_+|}{m} \qquad \mu_- = \dfrac{|S_-|}{m} \qquad \mu_0 = \dfrac{|S_0|}{m}$

↓ fraction correct

↓ fraction mistake

↓ what do thy care

$$\mathcal{L}_S^{exp}(\alpha f) = \mu_+ e^{-\alpha} + \mu_- e^{\alpha} + \mu_0$$

$\mathcal{L}(\alpha f)$ is convex in $\alpha$ **Yeh! Yey! Ci ! |> !**

$$0 = \frac{d\mathcal{L}}{d\alpha} = -\mu_+ e^{-\alpha} + \mu_- e^{\alpha} \quad \Longrightarrow \quad \alpha^* = \frac{1}{2} \log\left(\frac{\mu_+ m}{\mu_- m}\right)$$

Nobody is perfect assumption: $\mu_+ > 0 \quad \mu_- > 0$

# Further Insights

✱ If $\mu_+ = \mu_-$ then $\alpha = 0$ $\longrightarrow$ $f(x)$ is no better than a random predictor

✱ If $\mu_- > \mu_+$ then $\alpha < 0$ $\longrightarrow$ Negate the prediction of $f$

---

## Total Loss of $\alpha f(x)$ :

$$\mathcal{L}_S^{exp}(\alpha f) = \mu_+ \, e^{-\frac{1}{2} \log\left(\frac{\mu_+}{\mu_-}\right)} + \mu_- \, e^{\frac{1}{2} \log\left(\frac{\mu_+}{\mu_-}\right)} + \mu_0$$

$$= \mu_+ \sqrt{\frac{\mu_-}{\mu_+}} + \mu_- \sqrt{\frac{\mu_-}{\mu_+}} + \mu_0$$

$$= 2\sqrt{\mu_+ \mu_-} + \left(1 - (\mu_- + \mu_-)\right)$$

$e^{\frac{1}{2}\log a} = e^{\log \sqrt{a^2}} = \sqrt{a^2}$

$\mu_+ + \mu_- + \mu_0 = 1$

# Total loss of $\alpha f$

$$\mathcal{L}_S^{exp}(\alpha f) = \mu_+ \, e^{-\frac{1}{2} \log\left(\frac{\mu_+}{\mu_-}\right)} + \mu_- \, e^{\frac{1}{2} \log\left(\frac{\mu_+}{\mu_-}\right)} + \mu_0$$

$$= \mu_+ \sqrt{\frac{\mu_-}{\mu_+}} + \mu_- \sqrt{\frac{\mu_+}{\mu_-}} + \mu_0$$

$$= 2\sqrt{\mu_+ \mu_-} + \left(1 - (\mu_+ + \mu_-)\right)$$

$e^{\frac{1}{2} \log a} = e^{\log \sqrt{a}} = \sqrt{a}$

$\mu_+ + \mu_- + \mu_0 = 1$

No normalization $\Longrightarrow$ $m - \mu_+ - \mu_-$

$$\mathscr{L}(\alpha f) = 2\sqrt{\mu_+ \mu_-} + (1 - (\mu_+ + \mu_-))$$

$$= 1 - (\mu_+ - 2\sqrt{\mu^{+1}}\sqrt{\mu^{-1}} + \mu^-) \qquad (\sqrt{\mu^{-1}})^2$$

$$= 1 - (\sqrt{\mu^{+1}} - \sqrt{\mu^{-1}})^2$$

$\Delta Loss \equiv$ loss before using $f$ and after

$$\Delta \mathscr{L} \equiv \mathscr{L}(0) - \mathscr{L}(\alpha f) = (\sqrt{\mu^{+1}} - \sqrt{\mu^{-1}})^2$$

Reduction in loss has a closed form

Neat!

# Beyond a single predictor

1) Assume we calibrated $\alpha_1 h_1(x)$

2) Provided with a second predictor $h_2(x)$

3) Task: **combine** $\underbrace{\alpha_1 h_1(x)}_{\text{known}} + \alpha_2 h_2(x) \equiv f_2(x)$

     provided

$$\equiv f_1(x)$$

$$f_1 : x \to ? \{-\alpha_1, 0, \alpha_1\} \quad ?$$

4) General case: given $f_t(x) = \sum_{c=1}^{t} \alpha_c \cdot h_c(x)$ & $h_{t+1}(x)$

$$\text{find } \alpha_{t+1} \implies f_{t+1}(x) = \sum_{c=1}^{t+1} \alpha_c \cdot h_c(x)$$

$$= f_t(x) + \alpha_{t+1} h_{t+1}(x)$$

# Inductive Calibration

## As before define:

$$S_+^{t+1} = \{ i \mid h_{t+1}(x_i) = y_i \}$$

$$S_-^{t+1} = \{ i \mid h_{t+1}(x_i) = -y_i \}$$

$$S_0^{t+1} = \{ i \mid h_{t+1}(x_i) = 0 \}$$

## Generalize:

$$\mu_+ = \frac{1}{m} \sum_{i \in S_+} e^{-y_i f_t(x_i)}$$

$$\mu_- = \frac{1}{m} \sum_{i \in S_-} e^{-y_i f_t(x_i)}$$

$$\mu_0 = \frac{1}{m} \sum_{i \in S_0} e^{-y_i f_t(x_i)}$$

$$\mathcal{L}\left(f_{t+1}(x)\right) = \mu_+ e^{-\alpha_{t+1}} + \mu_- e^{+\alpha_{t+1}} + \mu_0$$

$$f_{t+1}^{(\cdot)} = f_t^{(\cdot)} + \alpha_{t+1} h_{t+1}^{(\cdot)}$$

Solution has the same form !

$$\Delta \mathcal{L}_{t+1} \equiv \mathcal{L}(f_t) - \mathcal{L}(f_{t+1}) = \left(\sqrt{\mu^+} - \sqrt{\mu^-}\right)^2$$

$$\left(\mu_+ + \mu_- + \mu_0\right) - \left(2\sqrt{\mu^+ \mu^-} + \mu_0\right)$$

# Importance Weights

* Given $f_t(x)$ define:

$$q_i^t \sim e^{-y_i f_t(x_i)}$$

  (importance weight:
  how "difficult" example $(x_i, y_i)$ is

  * Often $\{q_i^t\}$ is normalized $\sum_{i=1}^{m} q_i^t = 1$

  * Does not change analysis and algorithm

  * Simply implies $\mu_+^t + \mu_-^t + \mu_0^t = 1$

  Which can be obtained by simple scalinging

# Detour - SGD w/ Weights

$$\mathcal{L}(\omega) = \sum_{i=1}^{m} q_i \, f_i(\omega) \quad \text{s.t.} \quad \sum q_i = 1 \, ; \quad q_i > 0$$

Reduction : define $\tilde{f}_i(\omega) = q_i \, f_i(\omega)$

use SGD on : $\dfrac{1}{m} \sum_{i=1}^{m} \tilde{f}_i(\omega)$

## OR

Importance Sampling :   batch size $b$

Sample $b$ times such that example $i$ is chosen w/ probability $q_i$

# Boosting

Initialize: $\forall i \quad q_i^1 = \frac{1}{m} \quad f_0(x) \equiv 0$

For $t = 1, \ldots, T$:

Find $h_t(x)$ s.t. $\sum_i q_i^t y_i h_t(x_i)$ is large

Partition: $S \mapsto S_+^t, S_-^t, S_0^t$

Calculate: $\alpha_t = \frac{1}{2} \log\left(\frac{\mu_+^t}{\mu_-^t}\right)$

$\mu_+^t$
$\mu_-^t$
$\cancel{\mu_0^t}$

Update:

1. $f_{t+1}(x) = f_t(x) + \alpha_t h_t(x)$

2. $q_i^{t+1} = q_i^t e^{-y_i \alpha_t h_t(x_i)}$

$q_i^{t+1} = \begin{cases} q_i^t & h_t(x_i) = 0 \\ q_i^t e^{-\alpha} & h_t(x_i) = y_i \\ q_i^t e^{\alpha} & h_t(x_i) \neq y_i \end{cases}$

Equivalent to previous lecture

# Generalization

Instead of $h_t : \mathcal{X} \to \{-1, 0, +1\}$ can use

$$h_t(x) \in [-1, +1]$$   ⟹   Limit the power of base predictors

where as before

$$|h_t(x)| \quad \text{confidence}$$

$$\text{sign}(h_t(x)) \quad \text{predicted outcome}$$

$$S_+^t \triangleq \{ i \mid y_i \, h_t(x_i) > 0 \}$$

$$S_0^t \triangleq \{ i \mid y_i \, h_t(x_i) = 0 \}$$

$$S_-^t \triangleq \{ i \mid y_i \, h_t(x_i) < 0 \}$$

No further changes are required

Can replace exp-loss with log-loss (logistic-loss)

$$\mathcal{L}_S^{log}(f) = \frac{1}{m} \sum_{i=1}^{m} log(1 + e^{-y_i \cdot f(x_i)})$$

$$q_i^{t+1} \sim e^{-y_i \cdot f_{t+1}(i')} \sim q_i^t e^{-y_i \cdot \alpha_t h_t(x_i)}$$

$$q_i^{t+1} = \frac{1}{1 + exp(y_i f_t(x_i))} = \frac{q_i^t}{q_i^t + (1 - q_i^t) e^{y_i \alpha_t h_t(x_i)}}$$

Alternatively: $z_i^{t+1} \triangleq y_i f_{t+1}(x_i) = z_i^t + y_i \alpha_t h_t(x_i)$

$$q_i^{t+1} = \frac{e^{-z_i^{t+1}}}{Z} \qquad \text{Exp Loss} \qquad q_i^{t+1} = \left(1 + exp(z_i^{t+1})\right)^{-1} \qquad \text{Log Loss}$$

$$\text{Boosting}$$

Log-loss
Confidence-rated
W. H.

Initialize: $\forall i \quad q_i^1 = \frac{\cancel{k}}{\cancel{m}} \quad f_0(x) \equiv 0$

$\frac{1}{2}$

For $t = 1, \ldots, T$:

Find $h_t(x)$ s.t. $\sum_i q_i^t y_i h_t(x_i)$ is large

Partition: $S \mapsto S_+^t, S_-^t, S_0^t$

Calculate: $\alpha_t = \frac{1}{2} \log\left(\frac{\mu_+^t}{\mu_-^t}\right)$

Update:
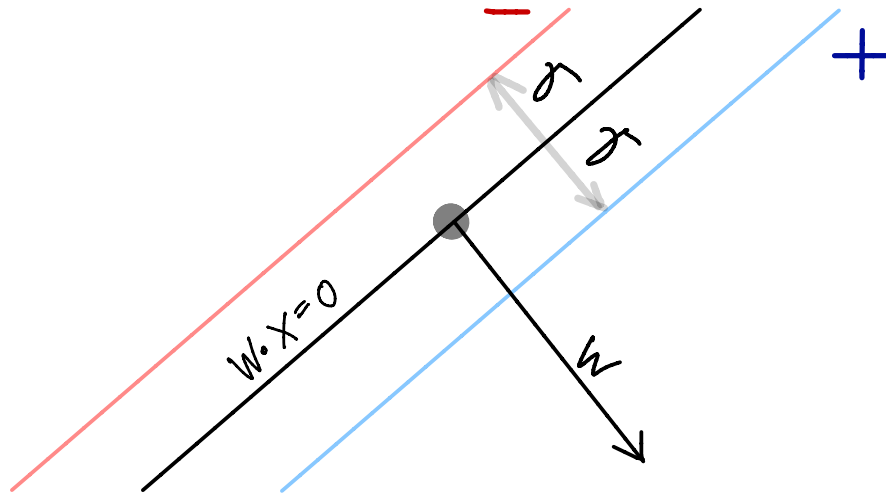
1. $f_{t+1}(x) = f_t(x) + \alpha_t h_t(x)$

2. $q_i^{t+1} = q_i^t e^{-y_i \alpha_t h_t(x_i)}$

2. $q_i^{t+1} = q_i^t \left(q_i^t + (1 - q_i^t) e^{y_i \alpha_t h_t(x_i)}\right)^{-1}$

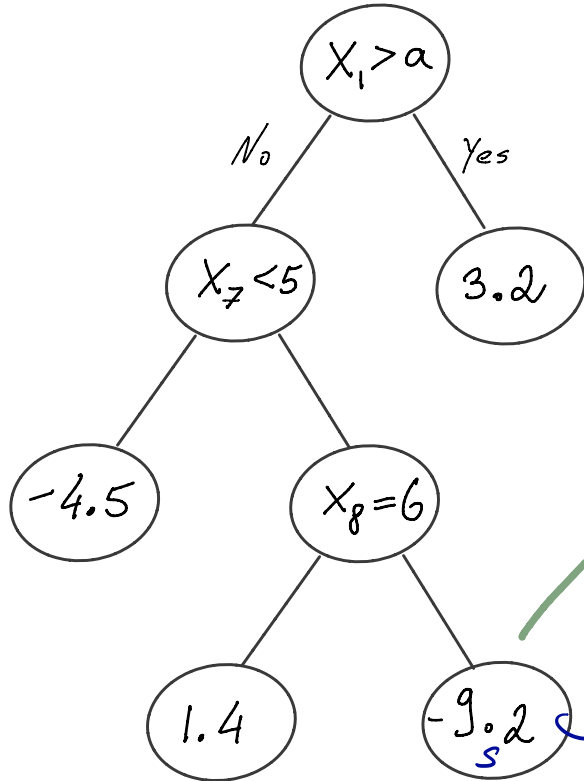# Detour - Linear Predictors w/ Abstantion

$$f(x) = w^T x \implies \tilde{f}(x) \text{ w/ abstantion \& continous}$$



$$W^T X = Z \qquad \tilde{f}(z) = \text{sign}(z)\left[|z| - \gamma\right]_+ = \begin{cases} z - \gamma & z > \gamma \\ 0 & |z| \leq \gamma \\ z + \gamma & z < -\gamma \end{cases}$$
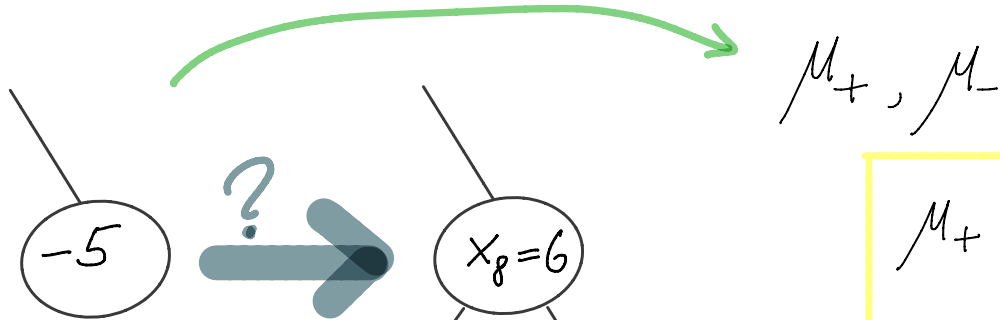
# Decision Trees

At each leaf there exists
a confidence-rated predictor



$$h_t(x) = \begin{cases} -9.2 & \text{if} \\ & x_1 \leq a \wedge x_7 \geq 5 \wedge x_8 = 6 \\ 0 & \text{O.W.} \end{cases}$$

$$\{i : y_i = +1\} = \mu_s^+$$

$$\{i : y_i = -1\} = \mu_s^-$$

$$\frac{1}{2} \log\left(\frac{\mu_s^+}{\mu_s^-}\right)$$

# Growing Decision Trees



$\mu_+ , \mu_-$

$$\mu_+ = \mu_+^L + \mu_+^R$$

$$\mu_- = \mu_-^L + \mu_-^R$$

$$\mu_0^R = \mu_+^L + \mu_-^L$$

$$\mu_0^L = \mu_+^R + \mu_-^R$$

$\mu_+^L , \mu_-^L , \mu_0^L$

$\mu_+^R , \mu_-^R , \mu_0^R$

Improvement in surrogate loss?

# Growing ~~Pain~~ Gain
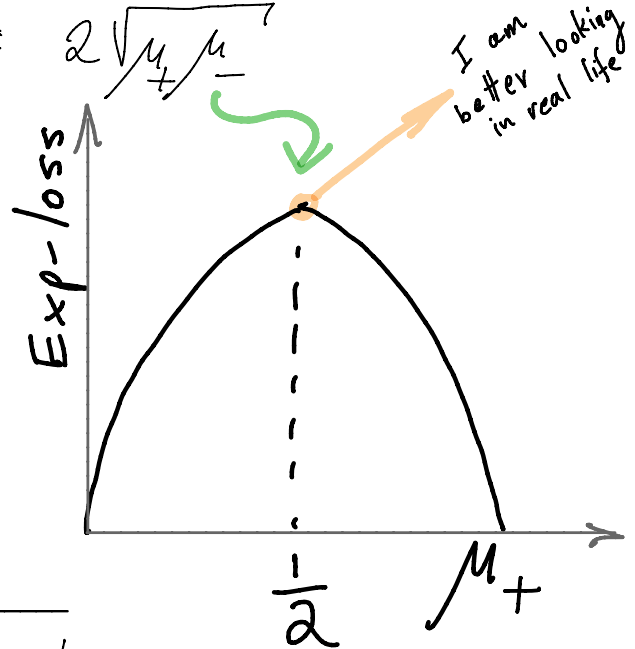
Assume $\mu_+ + \mu_- = 1$ for parent node (normalization)

Exp-loss for parent node: $2\sqrt{\mu_+ \mu_-}$

Exp-loss for children:

$$2\left(\sqrt{\mu_+^L \mu_-^L} + \sqrt{\mu_+^R \mu_-^R}\right)$$

Gain $\sim$

$$\sqrt{(\mu_+^R + \mu_+^L)(\mu_-^R + \mu_-^L)} - \sqrt{\mu_+^R \mu_-^R} - \sqrt{\mu_+^L \mu_-^L}$$
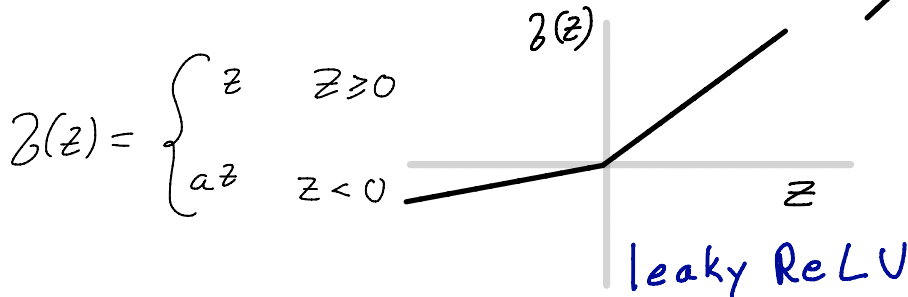
I am better looking in real life
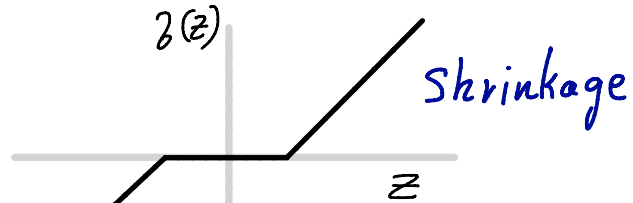
# Neural Networks

$$f_t(x) = \sum_{c=1}^{t} \alpha_i \, h_t(x) \implies \text{"Strong" hypothesis}$$

Suppose each $h_t(x)$ is of the form:

$$\mathcal{Z}(w^t \cdot x) \quad \text{where } \mathcal{Z}: \mathbb{R} \to \mathbb{R} \text{ is non-linear}$$

Examples:

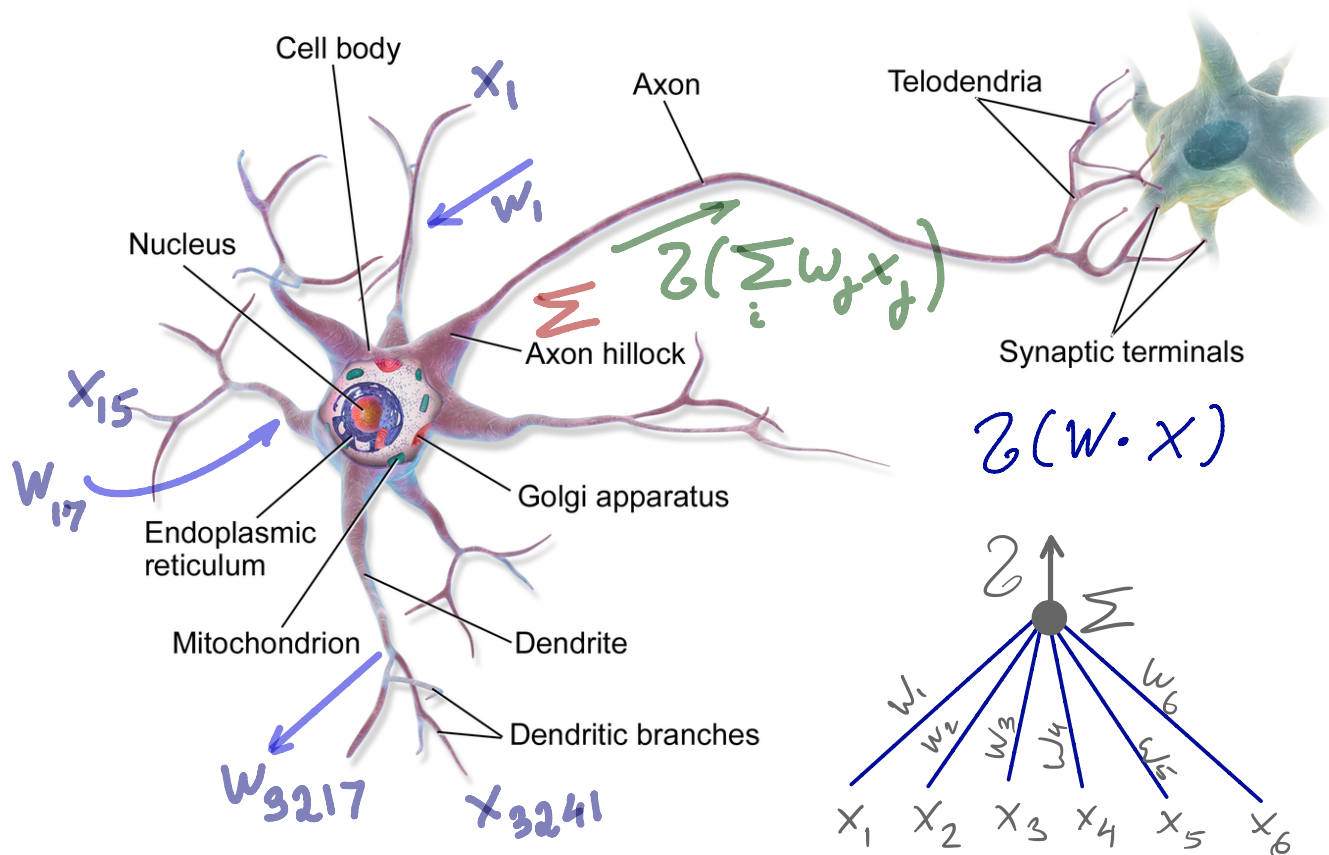$$\mathcal{Z}(z) = \begin{cases} z - \tau & z > \tau \\ 0 & |z| \leq \tau \\ z + \gamma & z < -\tau \end{cases}$$

$\mathcal{Z}(z)$

Shrinkage

$z$

$$\mathcal{Z}(z) = \begin{cases} z & z \geq 0 \\ az & z < 0 \end{cases}$$

$\mathcal{Z}(z)$

$z$

leaky ReLU

$$\mathcal{Z}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Sigmoid

"Neuron"

Cell body

$X_1$

$W_1$

Axon

Telodendria

Nucleus

$\Sigma$

$\mathcal{E}\left(\sum_i w_j x_j\right)$

Axon hillock

$X_{15}$

$W_{17}$

Synaptic terminals

$\mathcal{E}(W \cdot X)$

Golgi apparatus

Endoplasmic reticulum

Mitochondrion

Dendrite

$\mathcal{E}$

$\Sigma$

Dendritic branches

$W_{9217}$

$X_{3241}$

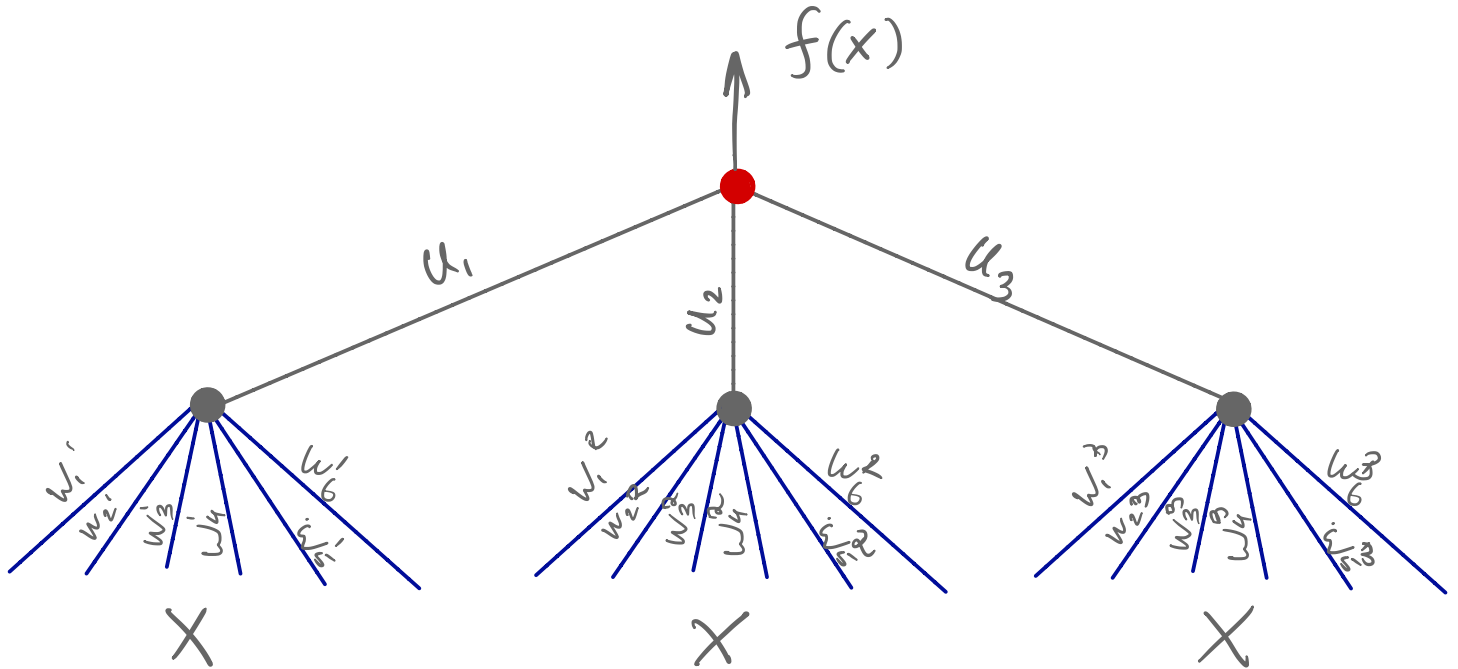$W_1$ $W_2$ $W_3$ $W_4$ $W_5$ $W_6$

$X_1$ $X_2$ $X_3$ $X_4$ $X_5$ $X_6$

# Multi layer perceptron (NN)

$$f(x) = \mathcal{Z}\left( u_1\, \mathcal{Z}(w^1 \cdot x) + u_2\, \mathcal{Z}(w^2 \cdot x) + u_3\, \mathcal{Z}(w^3 \cdot x)\right)$$

- You are likely to have many questions at this point...

- A few will be addressed @ IML

- Many will be left unanswered

- Since we don't know... the answer