# ① Adversarial Examples
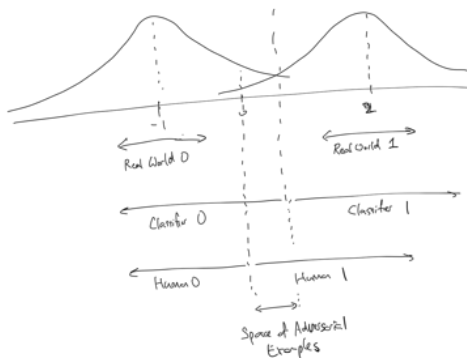
E.g. Logistic Regression Classifier on 1D Dataset ($x \in \mathbb{R}$)

BCE Loss $= -\frac{1}{m}\sum_{i=1}^{m} y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$

where $\hat{y} = \sigma(\omega x + b)$ where $\omega \in \mathbb{R}$, $b \in \mathbb{R}$.

Assume:
- $x | y=0 \sim N(-1, 1)$
- $x | y=1 \sim N(2, 1)$
- To humans, if $x < 0$, looks like $y=0$
  if $x \geq 0$, looks like $y=1$
  (This is a reasonable assumption. Data in the real world might lie in say
  $(-1.5, -0.5)$ for class $0$ and $(1, 3)$ for class $1$ most of the time and
  so $sgn(x)$ might be a good heuristic for humans!)



Real World 0     Real World 1

Classifier 0     Classifier 1

Human 0     Human 1

Space of Adversarial Examples

---

(A) Train Logistic Regression Classifier

After training, to find boundary, we need $P(y=0|x) = P(y=1|x) = 0.5$

$0.5 = \sigma(\omega x + b)$

$0.5 = \frac{1}{1 + e^{-\omega x - b}}$

$2 = 1 + e^{-\omega x - b}$

$e^{-\omega x - b} = 1$

$-\omega x - b = 0$

$x = -\frac{b}{\omega}$   #

---

(B) Naive Adversarial Example

Let's start with an example where $y=1$.

e.g. $x = 2$.

Now, we train with the same BCE loss but with label $y=0$.

Use BCE loss to backprop to $x$:

$x = x - \alpha \frac{\partial L}{\partial x}$

End Result (see Jupyter): Gets classified as $y=0$ ✓
but no longer "looks like" $y=1$ ✗

---

(C) Better Adversarial Example

Start with same example ($x = 2$)

Train with label $y = 0$

Loss $= $ BCE loss $+ \lambda (x_{new} - 2)^2$

End Result (see Jupyter): Gets classified as $y=0$ ✓
and still "looks like" $y=1$ ✓

---

(D) Black Box vs. White Box Attacks

White Box: Access to architecture and weights
   ↳ Like example above
   ↳ Can compute $\frac{\partial L}{\partial x}$.

Black Box: No access to architecture or weights

Black Box: No access ...

$\hookrightarrow$ Strategy: Generate a dataset $D = [(x_i, y_i)]_{i=1}^n$
where each datapoint $(x^{(i)}, y^{(i)})$ is taken by
sampling $x$ from some distribution and getting the
corresponding $y$ from the model $M$ that
you're trying to attack.

$\hookrightarrow$ Train new classifier $\hat{M}$ on dataset $D$. (like part (A))

$\hookrightarrow$ Do white box attack on $\hat{M}$ (like part (C))

$\hookrightarrow$ Hope the resulting $x$ is also an adversarial example for $M$.
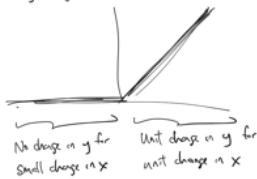
---

(E) Targeted vs. Non-targeted Attacks

Targeted: Want $M$ to misclassify to specific class $j$

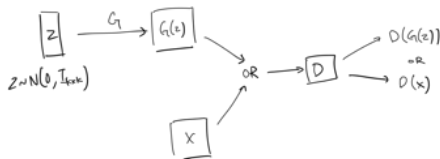Non-targeted: Want $M$ to misclassify to any class other than correct class

---

(F) Fast Gradient Sign Method

- (C) may be slow! $\Rightarrow$ Many gradient steps

- Idea: Just do $x = x - \varepsilon \, \text{sign}\left(\frac{\partial L}{\partial x}\right)$ where $L$ is original BCE loss.
  $\underset{\hookrightarrow}{\overline{\phantom{xx}}}$ Minimize perceptual difference
  $\hookrightarrow$ e.g. if $\left(\frac{\partial L}{\partial x}\right)_j$ is big in magnitude,
  without taking the sign, $x$ may no longer
  look like original $x$.

- Especially effective since most non-linear activations are usually
  in the linear regime (for fast training).

  e.g. $y = \text{ReLU}(wx+b)$



No change in $y$ for        Unit change in $y$ for
small change in $x$         unit change in $x$

---

② GANs



$z \sim N(0, I_{k\times k})$

Notes:
- $y$ is 1 if input is $x$.
- $y$ is 0 if input is $G(z)$.
- $D(G(z))$ and $D(x)$ are probabilities with range $(0,1)$.

---

$D$ is a classifier. $\Rightarrow$ Use BCE loss!

$$J^{(D)} = -\frac{1}{M} \sum_{i=1}^{m} y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$$

$$= -\frac{1}{M_{real}} \sum_{i=1}^{M_{real}} \log(D(x^{(i)})) - \frac{1}{M_{gen}} \sum_{i=1}^{M_{gen}} \log(1 - D(G(z^{(i)})))$$

$G$ is trying to fool $D$.

So $J^{(G)} = -J^{(D)}$

$$= \frac{1}{M_{real}} \sum_{i=1}^{M_{real}} \log(D(x^{(i)})) + \frac{1}{M_{gen}} \sum_{i=1}^{M_{gen}} \log(1 - D(G(z^{(i)})))$$

$\longrightarrow 0$ ($\because$ no $G$ in this term. Constant!)

---

But this formulation of $J^{(G)}$ is a saturating loss. This is because initially, discriminator is
doing very well ($\because$ much easier to train classifier as compared to generator).
$\therefore$ Initially, $D(G(z)) \approx 0.5$. Very quickly, $D(G(z)) \approx 0$ because discriminator
trains fast.

$$\frac{\partial J^{(G)}}{\partial h} = \frac{1}{M_{gen}} \sum_{i=1}^{M_{gen}} \frac{-D(1-D)}{1-D} \approx 0 \qquad \searrow 1 \text{ (same thing)}$$

Since $D$ is
a discriminator
between 2
classes, we
can write $D$ as

From Chain Rule, $\dfrac{\partial J^{(G)}}{\partial G_\theta} = \dfrac{\partial J}{\partial D_{out}} \cdot \underbrace{\dfrac{\partial D_{out}}{\partial D_{in}}}_{0} \cdot \dfrac{\partial h_{out}}{\partial h_{out}} \cdot \dfrac{\partial G_\theta}{\partial G_\theta}$    $\sigma'(h)$ where $h$ is last layer before applying sigmoid.

$\approx 0$.

Instead of minimizing the likelihood of discriminator being correct, we Maximize the likelihood of discriminator being wrong.

$\therefore$ $\min J^{(G)} = \min \left[ \dfrac{1}{M_{gen}} \sum\limits_{i=1}^{M_{gen}} \log\left(1 - D\left(G(z^{(i)})\right)\right) \right]$

$= \max \left[ \dfrac{1}{M_{gen}} \sum\limits_{i=1}^{M_{gen}} \log\left(D\left(G(z^{(i)})\right)\right) \right] = \max J'^{(G)}$

Once again, initially $D(G(z)) \approx 0$. Then,

$\dfrac{\partial J'^{(G)}}{\partial h} = \dfrac{1}{M_{gen}} \sum\limits_{i=1}^{M_{gen}} \dfrac{\cancel{D}(1-D)}{\cancel{D}} \approx 1$

We call $J'^{(G)}$ a non-saturating loss for the generator.