

**Bayes' Rule**  $p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} = \frac{p(x|\theta)p(\theta)}{\int_{\theta} p(x|\theta)p(\theta)d\theta}$

**Gaussian**  $p(X|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(X-\mu)^2}{2\sigma^2}\right]$

$p(x|\mu, \Sigma) = \frac{1}{\sqrt{2\pi}^d} \frac{1}{\sqrt{|\Sigma|}} \exp\left[-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right]$

$\ln(p(x|\mu, \Sigma)) = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln|\Sigma| - \frac{1}{2}(y - \mu)^T \Sigma(y - \mu)$

**Expected value**  $E[X] = \int_X xp(x)dx = \sum_{x \in X} xp(x)$

$E[aX] = aE[X]; E[XY] \stackrel{\text{indep.}}{=} E[X]E[Y]$

$E[X + Y] = E[X] + E[Y]$

**Variance**  $Var[X] = \int_x (x - \mu)^2 p(x)dx$

$Var[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$

**Norm**

$\|x\|_p := \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$

$\|x\|_2 := \sqrt{x_1^2 + \dots + x_n^2}$

$\|x\|_1 := |x_1| + \dots + |x_n|$

$\|x\|_{\infty} := \max_i |x_i|$

Efficient projections only exists for  $L_1, L_2$  and  $L_{\infty}$  norms.

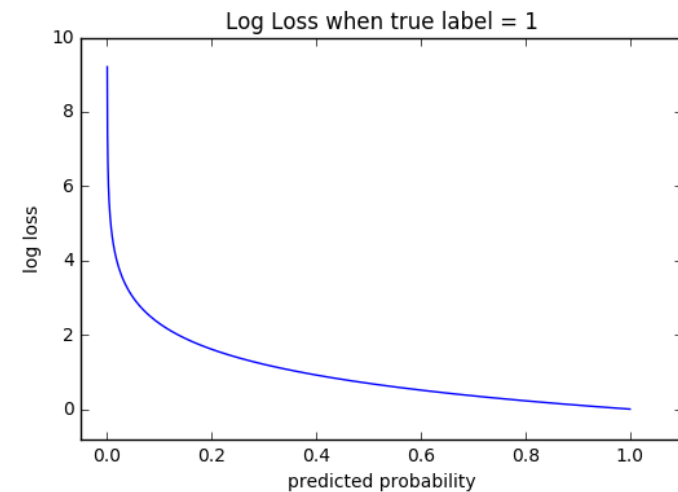
## 1 Deep Learning

### Loss

**Cross-entropy**  $\text{loss}(x, \text{true label}) = -\log\left(\frac{\exp(x[\text{true label}])}{\sum_j \exp(x[j])}\right)$

$= -\log(\text{softmax}(x[\text{true label}]))$

$= -x[\text{true label}] + \log\left(\sum_j \exp(x[j])\right)$ , where  $x$  is the *logit* output of the NN.



**NegativeLogLikelihoodLoss**  $\text{NLLLoss}(\text{logs}, \text{true label}) = -\text{logs}[\text{true label}]$ , where for logs the following should be used to make it equal to the *Cross-entropy loss*:

$\text{logs} = \log(\text{softmax}(x))$

## 2 Adversarial Examples

### Targeted FGSM (Fast Gradient Sign Method)

Intuition: Goal is to perturb the image such that the NN misclassifies the image to target label  $t$ . Therefore **reduce** the loss of the **target** label.

0. Target label  $t$ , true label  $s$ , generally  $t \neq s$

1. Compute perturbation:  $\eta = \epsilon \cdot \text{sign}(\nabla_x \text{loss}_t(x))$   
 $\nabla_x \text{loss}_t = \left( \frac{\partial \text{loss}_t}{\partial x_1}, \dots, \frac{\partial \text{loss}_t}{\partial x_n} \right)$ , where  $\text{loss}_t$  is the entry of the cross entropy loss vector for the target label and  $x$  is the image vector.
2. Perturb the input:  $x' = x - \eta$
3. Check if:  $f(x') = t$ , where  $f(x)$  is classification result of the NN for  $x$ .

### Untargeted FGSM (Fast Gradient Sign Method)

0. True label  $s$

1. Compute perturbation:  $\eta = \epsilon \cdot \text{sign}(\nabla_x \text{loss}_s(x))$   
 $\nabla_x \text{loss}_s = \left( \frac{\partial \text{loss}_s}{\partial x_1}, \dots, \frac{\partial \text{loss}_s}{\partial x_n} \right)$
2. Perturb the input:  $x' = x + \eta$
3. Check if:  $f(x') \neq s$

### PGD (Projected Gradient Descent)

Take  $k$  steps of **FGSM** each of size  $\epsilon$ . After each step project onto  $S(x)$ . By projecting, we mean that we find the closest point inside the  $S(x)$  ball (e.g.  $L_{\infty}$  ball). Here, closest is defined according to some norm (e.g.  $L_{\infty}$ ). In the region of  $S(x)$  we want all point to be classified with the **same** label.

$S(x) = \{x' | \|x - x'\|_{\infty} < \epsilon\}$

Note that the  $S(x)$  ball is of size  $\epsilon$  which is different than the  $\epsilon_{\text{step}}$  of the **FGSM step** and normally it holds that  $\epsilon_{\text{step}} < \epsilon$ . Projecting on the  $L_{\infty}$  ball is the same as clamping the values:  $x'_{\text{projected}} = \text{clamp}(x', \min = x - \epsilon, \max = x + \epsilon)$

Note that the resulting  $x'_{\text{projected}}$  can be inside the  $L_{\infty}$  ball.

### Minimization Problem for Defense:

find  $\theta$   
 minimize  $\rho(\theta)$   
 where  $\rho(\theta) = \mathbf{E}_{(x,y) \sim D} [\max_{x' \in S(x)} L(\theta, x', y)]$   
 in practice  $\rho(\theta) = \frac{1}{|D_a|} \sum_{(x,y) \in D_a} L(\theta, x, y)$

where  $D_a$  is dataset of adversarial examples

For which  $\rho(\theta)$  is the empirical risk (Loss) and the **outer** optimization (**min**) problem (Defense): *Find  $\theta$  that minimizes the high loss  $\rightarrow$  train robust classifier* (with normal SGD method  $\theta' = \theta - \epsilon_{\text{learning rate}} \cdot \nabla_{\theta} \rho(\theta)$ .)

Further,  $\max_{x' \in S(x)} L(\theta, x', y)$  the **inner** optimization (**max**)

problem (Attack): *Find adversarial  $x'$  achieves high loss  $\rightarrow$  adversarial attack.*

### Optimization Problem

Objective is to have a **small** perturbation  $\eta$ , such that the image is misclassified. Large perturbation is not wanted:

find  $\eta$   
 minimize  $\|\eta\|_p$   
 such that  $f(x + \eta) = t$   
 $x + \eta \in [0, 1]^n$

In general this is a hard problem to optimize with gradient descent. Therefore ease the constraints.

1. Use **objective function**:

$\text{obj}(x + \eta) \leq 0 \Rightarrow f(x + \eta) = t$

A correct objective function is a function that has  $\text{obj}(x') \leq 0 \iff p(t) \geq 0.5$

Sound objective functions for **2-class NN**:

$\text{obj}(x') = \text{loss}_t(x') - 1$

$\text{obj}(x') = \max(0, 0.5 - \text{softmax}(x')_t)$

For **k-class NN** this is:

$\text{obj}_k(x') = -\log_k(\text{softmax}(x')_t) - 1 = C \cdot \text{loss}_t(x') - 1$ , where  $C = \frac{1}{\log_2(k)}$  for cross-entropy loss with  $\log_2$ .

2. Replace norm with **proxy function** because the gradient of e.g. the inf norm is zero for most values except the maximum value.

Replace  $\|\eta\|_{\infty}$  with  $\sum_i \max(0, (\eta_i - \tau))$

If all entries are less than  $\tau$  then the entire expression is zero. Note: When  $\tau$  is large, the gradient is similar to the gradient of  $\|\eta\|_{\infty}$ . Start with large  $\tau$  and lower after each iteration.

3. Clamp perturbed image back to box domain after optimization.

Then the optimization becomes:

find  $\eta$   
 minimize  $\|\eta\|_p + c \cdot \text{obj}(x + \eta)$   
 such that  $x + \eta \in [0, 1]^n$

### Diffing Networks

Given two NN trained to learn same function. Perturb Input  $x$  such that  $\text{class}(f_1(x')) \neq \text{class}(f_2(x'))$

Use the following objective function, where  $f_i(x')_t$  is the softmax output of NN  $i$  w.r.t.

**while**  $\text{class}(f_1(x)) = \text{class}(f_2(x))$ :

$\text{obj}(x) = f_1(x)_t - f_2(x)_t \rightarrow$  Use as Loss

$x = x + \epsilon \cdot \frac{\partial \text{obj}(x)}{\partial x} \rightarrow$  **Maximize** Loss

**return**  $x$

### 3 Logic

Goal: Want to query NN such that some logical constraints are satisfied.

Problem: Formulating this as a constrained problem is hard to solve and times out for large NN.

Solution: Translate logical constrain into loss.

#### Translations

$\forall x$ , if  $T(\phi)(x) = 0 \Rightarrow \phi(x)$  satisfied,

where  $\phi(x)$  is logical formula

Use the following constrains to generate loss function.

Logical Term	Translation	Logic Negation ( $\neg$ )
$t_1 = t_2$	$ t_1 - t_2 $	$t_1 < t_2 \vee t_2 < t_1$
$t_1 \leq t_2$	$\max(0, t_1 - t_2)$	$t_1 > t_2$
$t_1 < t_2$	$T(t_1 + \epsilon \leq t_2)$	$t_1 \geq t_2$
$t_1 \neq t_2$	$T(t_1 < t_2 \vee t_2 < t_1)$	$t_1 = t_2$
$\varphi \vee \psi$	$T(\varphi) \cdot T(\psi)$	$\neg\varphi \wedge \neg\psi$
$\varphi \wedge \psi$	$T(\varphi) + T(\psi)$	$\neg\varphi \vee \neg\psi$

Problem: When dealing with **real** values in logical domain and **floats** in loss domain one has to assign  $\epsilon$  to smallest machine value. Thereafter Translation is only valid in one direction ( $T(\phi)(x) = 0 \Rightarrow \phi(x)$  satisfied). It no longer holds that when  $\phi(x)$  satisfied, that the loss is 0.

E.g.  $t_1 < t_2 \Rightarrow T(\phi) = \max(0, t_1 + \epsilon - t_2)$

Use,  $t_1 = t_2 - \frac{\epsilon}{2}$ . It holds that  $t_1 < t_2$  but Translation is not satisfied since  $T(\phi) = \max(0, t_2 - \frac{\epsilon}{2} - t_2) = \frac{\epsilon}{2} \neq 0$

Therefore,  $T(\phi)(x) = 0 \Rightarrow \phi(x)$  satisfied,

but  $\phi(x)$  satisfied  $\nRightarrow T(\phi)(x) = 0$

#### Example

Goal: Find an image  $i$  which gets classified to 9 where the image  $i$  is within some distance of the image deer.

0. Logical Formula:

$$\phi(i) = \bigwedge_{j=1, j \neq 9}^k \text{NN}(i)[j] < \text{NN}(i)[9] \wedge \|i - \text{deer}\|_{\infty} < 25 \wedge \|i - \text{deer}\|_{\infty} > 5$$

1. Translation into loss:

$$T(\phi) = \sum_{j=1, j \neq 9}^k \max(0, \text{NN}(i)[j] + \epsilon - \text{NN}(i)[9]) + \max(0, \|i - \text{deer}\|_{\infty} + \epsilon - 25) + \max(0, 5 + \epsilon - \|i - \text{deer}\|_{\infty})$$

2. Train Network with SGD

#### Train NN with Logic

0. Goal: Want to enforce property  $\phi$ . Find weights for NN, such that expected value of the property increases:

$$\begin{aligned} &\text{find } \theta \\ &\text{maximize } \rho(\theta) \\ &\text{where } \rho(\theta) = \mathbf{E}_{s \sim D} [\forall \mathbf{z} \cdot \phi(\mathbf{z}, s, \theta)] \end{aligned}$$

1. Translate into loss:

$$\begin{aligned} &\text{find } \theta \\ &\text{minimize } \rho(\theta) \\ &\text{where } \rho(\theta) = \mathbf{E}_{s \sim D} [T(\phi)(z_{\text{worst}}, s, \theta)] \\ &\text{and } z_{\text{worst}} = \arg \min_z (T(\neg\phi)(z, s, \theta)) \end{aligned}$$

Inner minimization: Find worst violation of property.

Outer minimization: Find weight such that worst violation is minimized.

Intuitively, we are trying to get the worst possible violation of the formula and then to find a network that minimizes its effect.

2. Solve inner minimization by splitting loss:

$$\begin{aligned} \text{e.g. } \text{loss}(z, x, \theta) &= \max(0, \|x - z\|_{\infty} - \epsilon) \\ &\quad + \max(0, \text{NN}_{\theta}(z)[3] - \delta) \end{aligned}$$

$\rightarrow$  split loss!

2.1. Solve with PGD:

$$\min_z \text{loss}(z, x, \theta) = \max(0, \text{NN}_{\theta}(z)[3] - \delta)$$

2.2. Project  $z$  back onto the  $L_{\infty}(x, \epsilon)$  ball

### 4 Certify AI - Abstract Domains

**Sound:** Correct Approximation of NN.

**Precise:** Approximation is superset of NN, but should not approximate too much, otherwise can not verify NN.

**Efficient:** Efficient to compute

#### Interval Domain

Input  $x$  is in the form of  $x = [a, b]$ .

#### Operation Rules

Addition:  $x_1 + x_2 = [x_1[0] + x_2[0], x_1[1] + x_2[1]]$

Subtract.:  $[x_1, x_2] - [y_1, y_2] = [x_1 - y_2, x_2 - y_1]$

Addition Scalar:  $x_1 + a = [x_1[0] + a, x_1[1] + a]$

Multipl.:  $[x_1, x_2] \cdot [y_1, y_2] = [\min(x_1 y_1, x_1 y_2, x_2 y_1, x_2 y_2), \max(x_1 y_1, x_1 y_2, x_2 y_1, x_2 y_2)]$

Multipl Scalar:  $x_1 \cdot a = [l, u]$ ,  
with  $l = \min(x[0] \cdot a, x[1] \cdot a)$ ,  $u = \max(x[0] \cdot a, x[1] \cdot a)$

Funct.:  $f([y_1, y_2]) = [\min\{f(y_1), f(y_2)\}, \max\{f(y_1), f(y_2)\}]$

Lower Equal:  $\leq ([l_1, u_1], [l_2, u_2]) = ([l_1, u_1] \sqcap_i [-\infty, u_2], [l_1, \infty] \sqcap_i [l_2, u_2])$

#### Zonotope Abstrac Domain

$$\hat{m} = a_0^m + \sum_{i=1}^k a_i^m \epsilon_i$$

$\epsilon_i$ : noise terms ranging  $[-1, 1]$  shared between abstract neurons

$a_i^n$ : real number that controls magnitude of noise  
Centered around  $a_0^m$

### Operation Rules

#### Multiplication with scalar

$$(a_0^n + \sum_{i=1}^k a_i^n \epsilon_i) \cdot C = (C \cdot a_0^n + \sum_{i=1}^k C \cdot a_i^n \epsilon_i), C \in \mathbb{R}$$

#### Multiplication of two variable

$$\begin{aligned} &\left(a_0^n + \sum_{i=1}^k a_i^n \epsilon_i\right) \cdot \left(a_0^m + \sum_{i=1}^k a_i^m \epsilon_i\right) = (a_0^n \cdot a_0^m) + \\ &\sum_{i=1}^k (a_i^n \cdot a_0^m + a_i^m \cdot a_0^n) \cdot \epsilon_i + \sum_{i=1}^k \sum_{j=1}^k a_i^m \cdot a_j^n * \epsilon_i \cdot \epsilon_j \end{aligned}$$

where  $\epsilon_i \cdot \epsilon_j$  becomes new variable  $\epsilon_{i,j}$  and

$$\epsilon_{i,j} \in [-1, 1] \text{ if } i \neq j$$

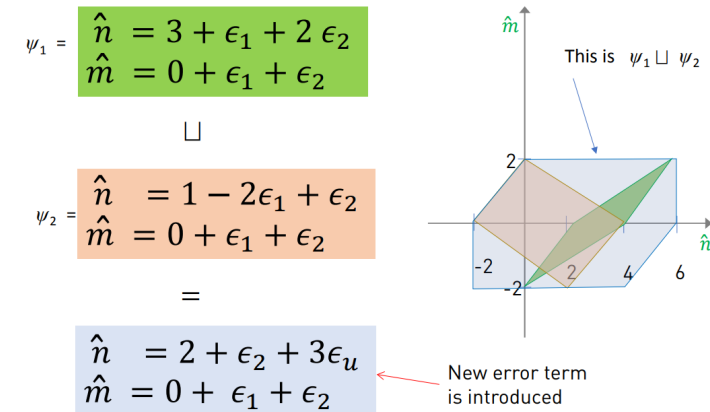
$$\epsilon_{i,j} \in [0, 1] \text{ if } i = j$$

#### Summation

$$\begin{aligned} &\left(a_0^n + \sum_{i=1}^k a_i^n \epsilon_i\right) + \left(a_0^m + \sum_{i=1}^k a_i^m \epsilon_i\right) = (a_0^n + a_0^m) \\ &\quad + \sum_{i=1}^k (a_i^n + a_i^m) \cdot \epsilon_i \end{aligned}$$

#### Join

Operation is not closed. Example of the Operation:



#### ReLU

$$f_{ReLU}^{\#} = \text{ReLU}_2^{\#}(b) \circ \text{ReLU}_1^{\#}(a) \text{ (Affine)}$$

$\text{ReLU}_i^{\#}(x_i)(\psi) = \psi_{\{x_i \geq 0\}} \sqcup \psi_{\{x_i < 0\}}$ ,  
where  $\psi_{\{x_i \geq 0\}} = (\psi \sqcap \{x_i \geq 0\})$ , for which  $\sqcap$  is not defined.

$$\text{and } \psi_{\{x_i < 0\}} = \begin{cases} [[x_i = 0]](\psi) & \text{if } (\psi \sqcap \{x_i < 0\}) \neq \perp \\ \perp & \text{otherwise} \end{cases}$$



```

1 def main(){ \\ Program with expectation of Pi.
2     x := uniform(0,1);
3     y := uniform(0,1);
4     within_radius := x*x+y*y <= 1;
5     return 4*within_radius; \\Prob. that point ...
6                             \\is within radius.
7 }

```

## Differential Privacy

### Epsilon - Differential Privacy

$$\frac{\Pr[F(x) \in \Phi]}{\Pr[F(x') \in \Phi]} \leq \exp(\varepsilon)$$

, where  $F(x)$  is the output of the database query including randomization,  $x'$  is a database that differs on a single element w.r.t. to  $x$  and  $\Phi$  is the secret.

### 7 Programming by Examples (PBE)

A new frontier in AI where one learns an interpretable program from user-provided examples.

Requires very few input-output examples. Assumes the given examples are representatives.

#### PBE problem definition

**Given:** A domain specific Language (DSL) & set of input-output examples.

**Goal:** Learn a function over the DSL which is consistent with the provided examples.