

Let n be the number of steps between evaluations.

Let p be the "patience," the number of times to observe worsening validation set error before giving up.

Let θ_0 be the initial parameters.

$\theta \leftarrow \theta_0$

$i \leftarrow 0$

$j \leftarrow 0$

$v \leftarrow \infty$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

while $j < p$ do

 Update θ by running the training algorithm for n steps.

$i \leftarrow i + n$

$v' \leftarrow \text{ValidationSetError}(\theta)$

 if $v' < v$ then

$j \leftarrow 0$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

$v \leftarrow v'$

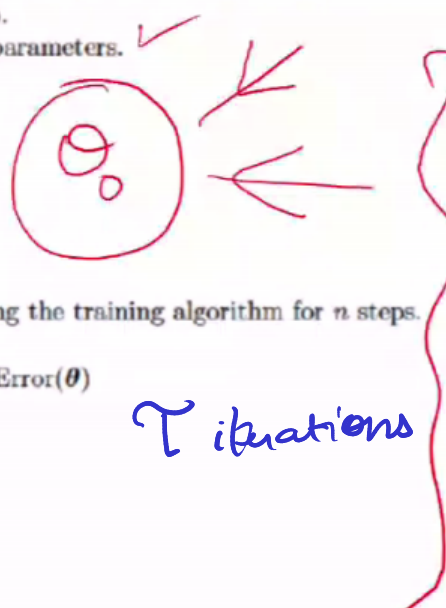
 else

$j \leftarrow j + 1$

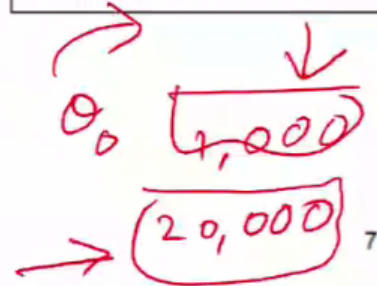
 end if

end while

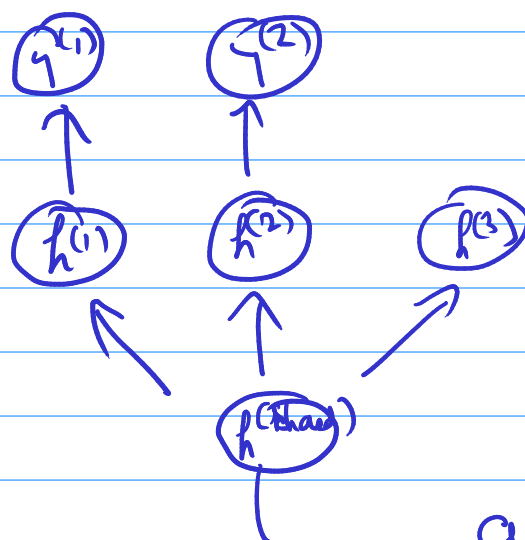
Best parameters are θ^* , best number of training steps is i^*



Algorithm determines the best amount of time to train. The meta algorithm is a general strategy that works well with a variety of training algorithms and ways of quantifying error on the validation set.



(Provides)
Note: SGD does implicit regularization



no. of steps
for early
stopping is
a hyperparameter

find this by training
a small subset, then
use the number to train whole model
(Shut up & TRAINNN)

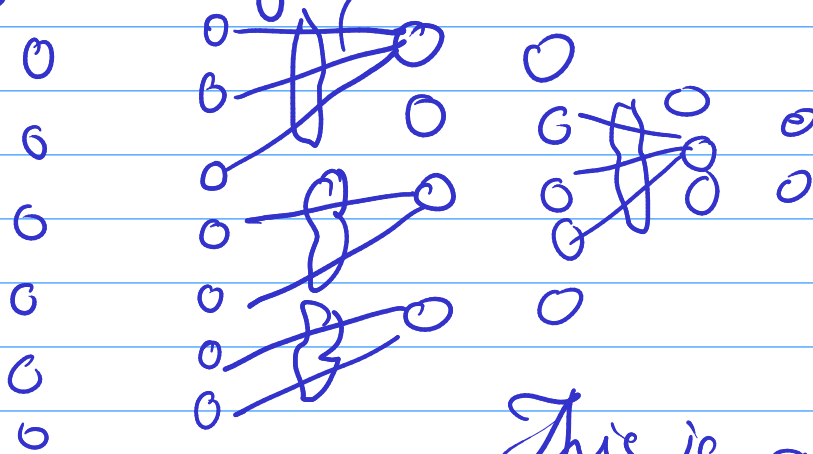
Flavor #2

Let $X^{(\text{train})}$ and $y^{(\text{train})}$ be the training set.
 Split $X^{(\text{train})}$ and $y^{(\text{train})}$ into $(X^{(\text{subtrain})}, X^{(\text{valid})})$ and $(y^{(\text{subtrain})}, y^{(\text{valid})})$ respectively.
 Run early stopping (algorithm 7.1) starting from random θ using $X^{(\text{subtrain})}$ and $y^{(\text{subtrain})}$ for training data and $X^{(\text{valid})}$ and $y^{(\text{valid})}$ for validation data. This updates θ .
 $\epsilon \leftarrow J(\theta, X^{(\text{subtrain})}, y^{(\text{subtrain})})$
 while $J(\theta, X^{(\text{valid})}, y^{(\text{valid})}) > \epsilon$ do
 Train on $X^{(\text{train})}$ and $y^{(\text{train})}$ for n steps.
 end while

θ^* for this
 don't throw θ^* away

70%

Weight sharing



This is also a regularizer,
 overfitting doesn't happen

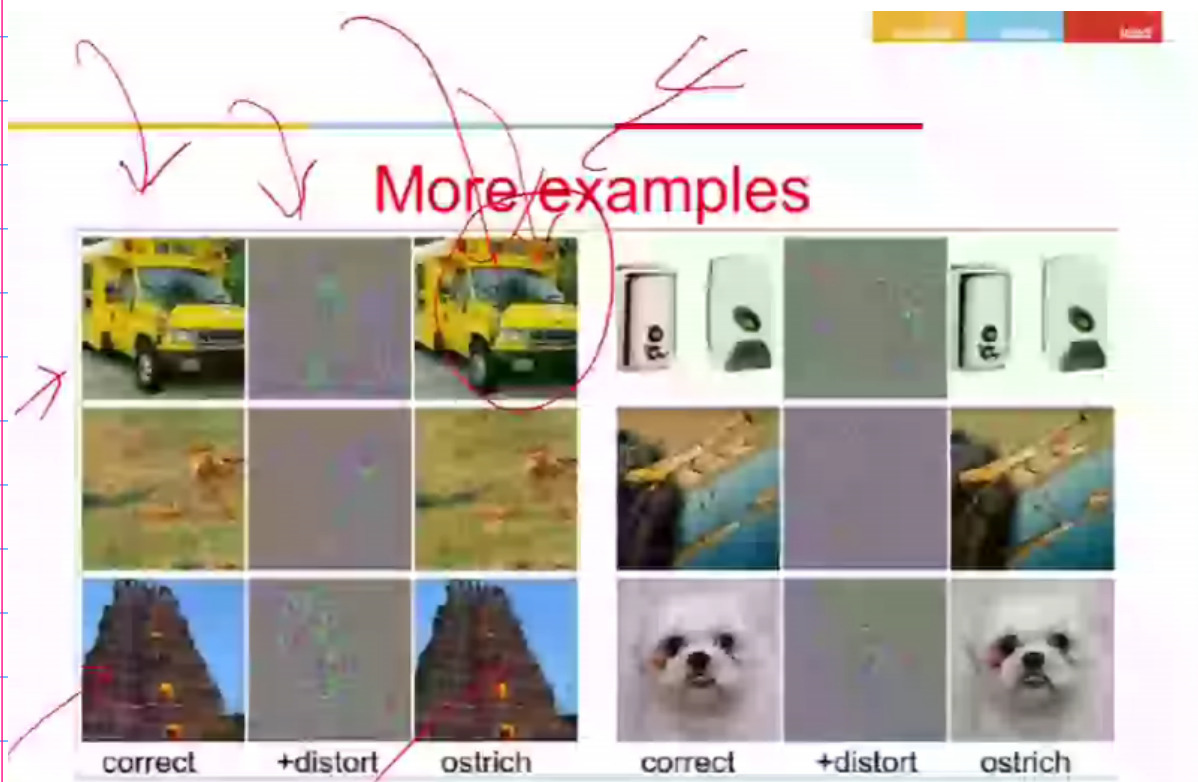
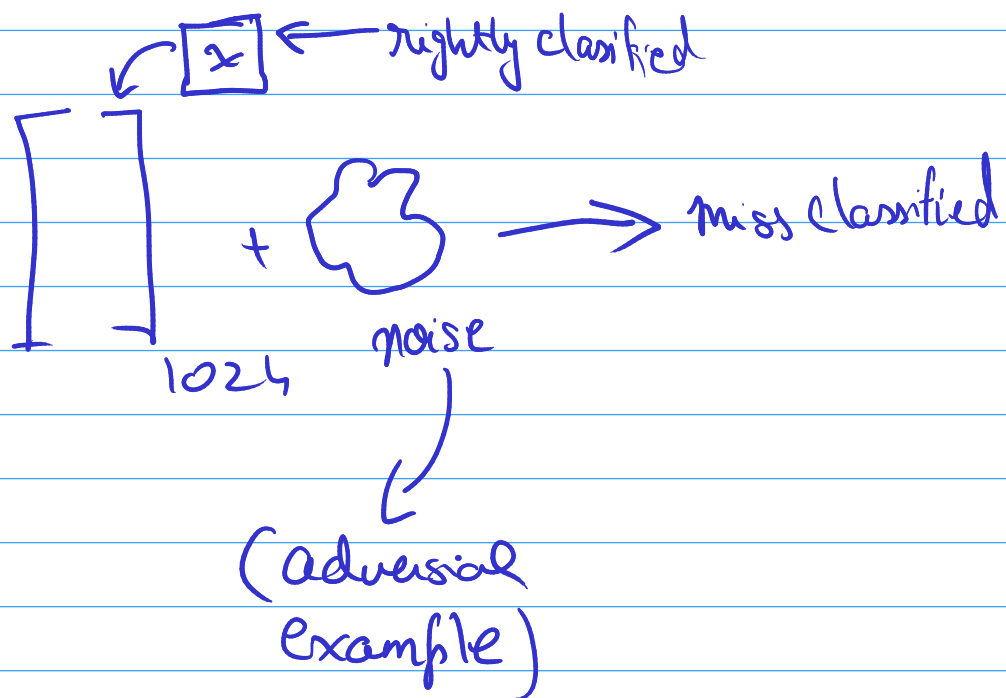
(adversarial)

Adversarial Training

Szegedy

2014

\times ← rightly classified



goodfellow 2014, explaining adversarial -- a composition of functions.
 but relu type functions composed can be incredibly linear most of the time

a heavily linear function will look like a lin. comb. of weights

a close neighbourhood can be $W.(x + \epsilon)$ (x, w and ϵ are vectors)

$$f(x + \epsilon) \approx \frac{1}{x + \epsilon} \approx \frac{1}{x}$$

if f is quadratic, ϵ will become ϵ^2 etc.
 and vanish

Energy layer

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_{30000} \end{bmatrix} \left(\begin{bmatrix} x_1 \\ \vdots \\ x_{30000} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_{30000} \end{bmatrix} \right)$$

ϵ is non linear

$\Rightarrow x$ becomes high non linear

$$f \left(\begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1100} \end{bmatrix} + \begin{bmatrix} \epsilon \\ \epsilon \\ \vdots \\ \epsilon \end{bmatrix} \right) \text{ changes quite a lot}$$

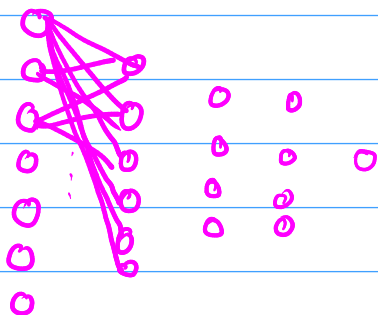
$$f(x) - f(x + \epsilon)$$

Apply 'f' on this

$$g = f_0 \circ f_1 \circ f_2 \dots f_n(x)$$

$$g(y) - g(y + \epsilon) \text{ will be way more than } f(x) - f(x + \epsilon)$$

$$f = \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_{100} x_{100}$$



$$f \approx \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_{100} x_{100}$$

Adding noise

$$\omega^T x + \omega^T \delta$$

$$\omega^T x + \omega^T \begin{bmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_{100} \end{bmatrix}$$

the quantity

let

$$\begin{aligned} \delta_i &\text{ be } +\epsilon && \text{if } \omega \geq 0 \\ \delta_i &\text{ be } -\epsilon && \text{if } \omega < 0 \end{aligned}$$

Whitebox
attack

$$x = x + \delta$$

$$\delta = \epsilon \begin{bmatrix} \text{sign}(\omega_1) \\ \vdots \\ \text{sign}(\omega_{100}) \end{bmatrix}$$

hyperbolic

$$\omega^T x + \epsilon \sum_{i=1}^{100} |\omega_i| = \omega^T x + \epsilon (|\omega_1| + |\omega_2| + \dots + |\omega_{100}|)$$

$$x \rightarrow (x + \delta)$$

$$f(x) \rightarrow f(x + \delta)$$

even

Since f is linear, if x is close to x'

$f(x')$ can be completely different

BUT Why?? (This seems contradictory)
you can't see this lol

Basically, this is an explanation, not a proof
for adversarial examples

It works for not all examples, f is normally piecewise

We are also NOT comparing with linear polynomial & exponential; we are just saying linearity itself makes values grow!
Compared to sigmoid etc.

SVMs are also linear but have large decision regions, so they don't face this problem.

⇒ Training on adversarial examples makes model more robust