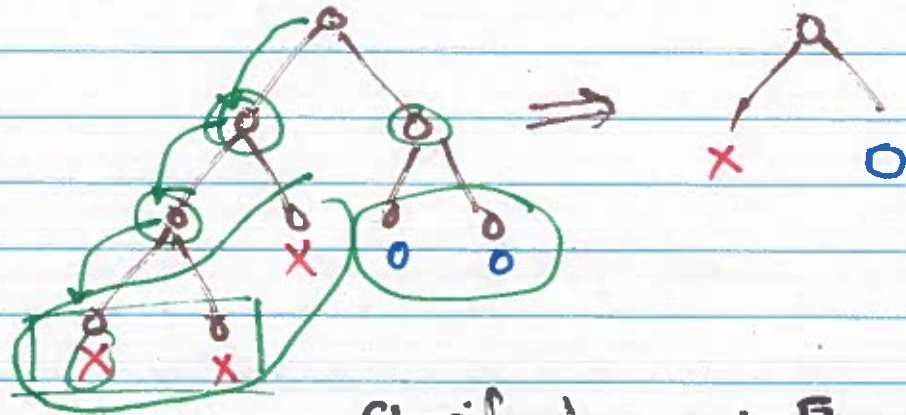


Approx K-NN: No backtrack



## Classification and Regression Tree CART

Goal: build a tree that's

- (1) as small as possible, and
- (2) only has 1 label per leaf node ("pure" leaves)

NP-Hard.

Idea: split recursively to minimize  
"impurity" at each node of tree.

heuristic that measures how "diverse" the labels  
are in the dataset.

## Gini Impurity. (Classification)

Empirical distribution over labels  $\mathcal{Y}$

$$p_k = \frac{|\{(x,y) \in S \mid y=k\}|}{|S|} = \frac{\text{\# of examples w/ label } k}{\text{\# of examples.}}$$

$$\text{Gini} \Rightarrow G(S) = \sum_k p_k (1-p_k)$$

\# of ways to draw 2 examples with unequal labels

$$\sum_k \sum_{l \neq k} \underbrace{|\{(x,y) \in S \mid y=k\}| \cdot |\{(x,y) \in S \mid y=l\}|}_{\substack{\text{\# with label } k \\ \text{\# with label } l}}$$

$$= \sum_k (\text{\# with label } k) \cdot \sum_{l \neq k} (\text{\# with label } l)$$

$$= \sum_k (\text{\# with label } k) \cdot (n - (\text{\# with label } k))$$

$$= \sum_k (n p_k) \cdot (n - (n \cdot p_k)) = n^2 \cdot \sum_k p_k (1-p_k)$$

$$= \sum_k (n p_k) \cdot (n - (n \cdot p_k)) = n^2 \cdot \sum_k p_k (1-p_k)$$

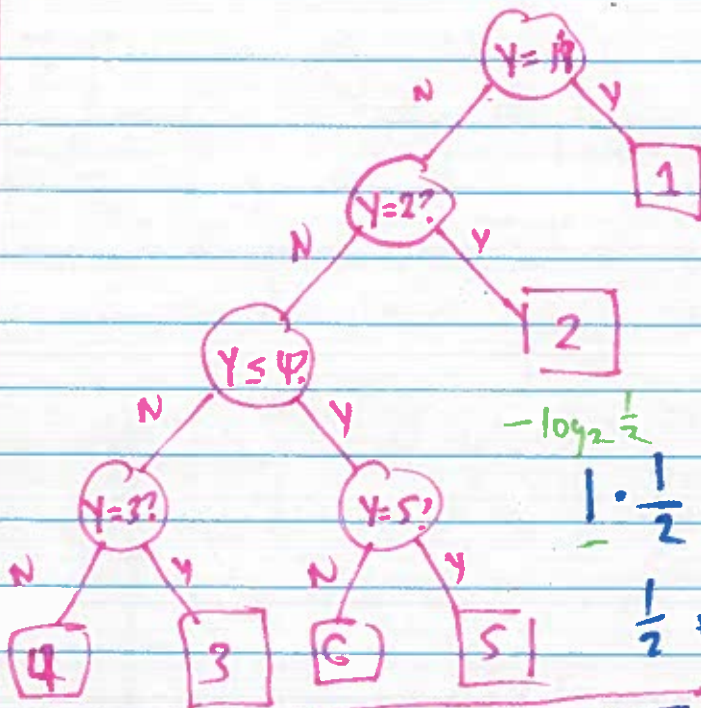
$$G^T(S) = \frac{|S_L|}{|S|} G^T(S_L) + \frac{|S_R|}{|S|} G^T(S_R)$$

$\wedge$   
 $S_L \quad S_R$



# Entropy.

$$P_1 = \frac{1}{2} \quad P_2 = \frac{1}{4}, \quad P_3 = P_4 = P_5 = P_6 = \frac{1}{16}$$



$$-\log_2 \frac{1}{2} \quad -\log_2 \frac{1}{4} \quad -\log_2 \frac{1}{16}$$

$$1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{4} + 4 \cdot \left(\frac{4}{16}\right)$$

$$\frac{1}{2} + \frac{1}{2} + 1 = 2$$

$$H(S) = \sum_k -P_k \log_2 P_k$$

$H(S)$  is a ~~convex~~  $D_{KL}(P \parallel \text{uniform})$   
 monotonic function of  $D_{KL}(P \parallel q)$

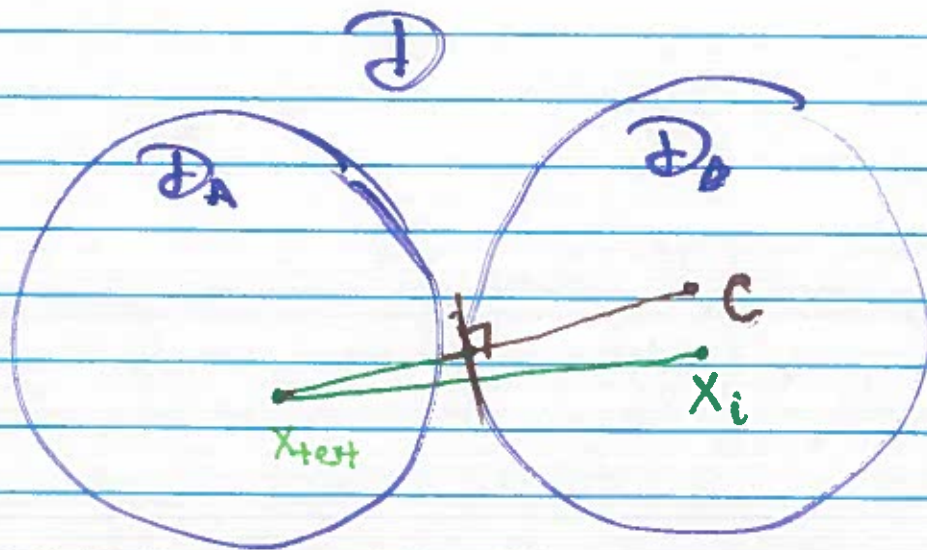
$$H^T(S) = \frac{|S_L|}{|S|} \cdot H^T(S_L) + \frac{|S_R|}{|S|} \cdot H^T(S_R)$$

$$D_{KL}(P \parallel q) = \sum_k P_k \log \left( \frac{P_k}{q_k} \right)$$

H

$$D_{KL}(P \parallel \text{uniform}) = \sum_k P_k \log \left( \frac{1/c}{P_k} \right) = \sum_k P_k \log \frac{1}{c} - \sum_k P_k \log P_k$$

# Ball Trees





## Impurity

- Gini impurity  $\approx$  how often will random labels disagree
- Entropy  $\approx$  information
- Least squares

## ID3

Input: Dataset  $\mathcal{D}$

Search all features:

search all splits for feature:

evaluate impurity (entropy)

pick feature/split that minimizes impurity

construct node

recursively call ID3 on subsets

## How to split?

- categorical features  $\rightarrow$  1 child per category
- real-valued features  $\rightarrow$  threshold

Base cases: if  $|\mathcal{D}| = 1$

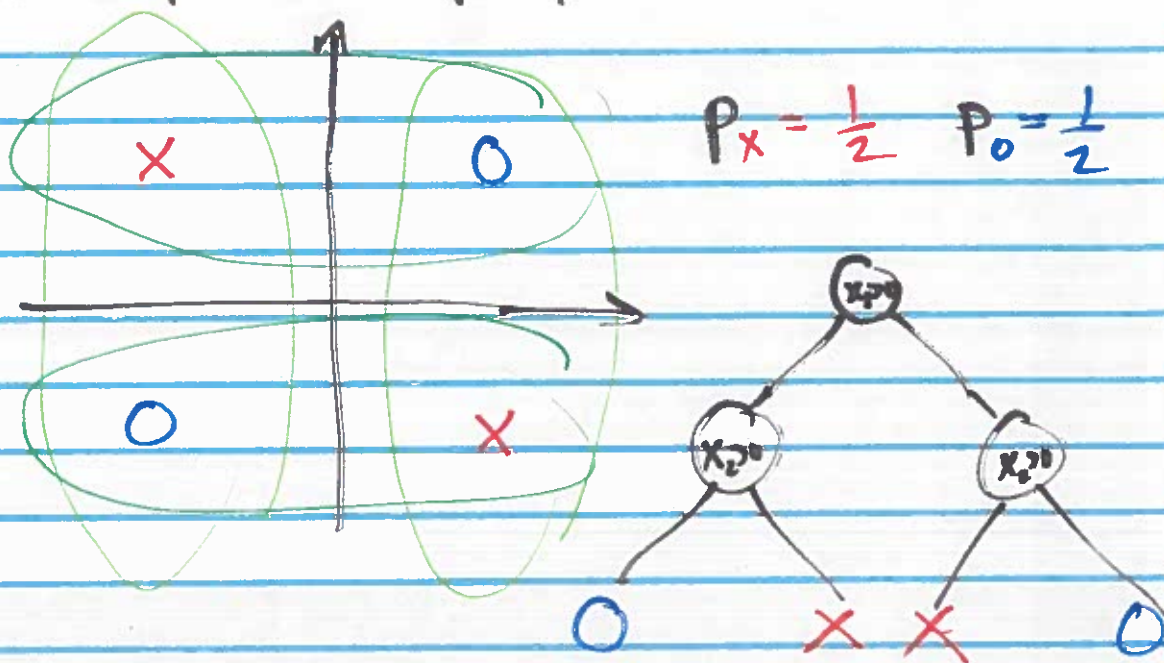
classification  $\rightarrow$  majority  
 $\rightarrow$  regression  $\rightarrow$  average

1.  $\exists$  exists some  $\hat{x}$  s.t.  $\forall (x,y) \in \mathcal{D}, x = \hat{x}$

2. exists  $\hat{y}$  s.t.  $\forall (x,y) \in \mathcal{D}, y = \hat{y} \Rightarrow$  leaf  $\hat{y}$

3.  $|\mathcal{D}| = 0$ , predict majority/average of "parent" dataset

Why not stop when impurity doesn't decrease.



Decision tree inference time is  
proportional to depth  
↓  
inference very fast!

Overfitting:

↑  
bias<sup>2</sup> + variance + noise



Ensembling: average the prediction of some models

- draw  $m$  independent datasets  $D_1, \dots, D_m$
- for each dataset:
  - run ID3  $\rightarrow$  hypothesis  $h_i$

- output  $\hat{h}(x) = \frac{1}{m} \sum_{i=1}^m h_i(x)$

How does this effect: bias<sup>2</sup> variance noise?  
same  $\downarrow$  factor of  $m$  same

$$\begin{aligned}\text{Var}\left(\frac{1}{m} \sum_{i=1}^m h_i(x)\right) &= \frac{1}{m^2} \text{Var}\left(\sum_{i=1}^m h_i(x)\right) \\ &= \frac{1}{m^2} \sum_{i=1}^m \text{Var}(h_i(x)) \\ &= \frac{1}{m^2} \cdot m \cdot \text{Var}(h_i(x)) \\ &= \frac{1}{m} \text{Var}(h_i(x))\end{aligned}$$

## Bootstrap Aggregating

instead of drawing from source dist,  
we draw with replacement from  $\mathcal{D}$

given  $\mathcal{D}$  ~~sampled~~  
draw  $m$  ~~int~~ datasets

draw  $m$  datasets from  $\mathcal{D}$  each of size  $n$   
for each,  $\mathcal{F}$  train a decision tree (ID3)

Average:  $h(x) = \frac{1}{m} \sum_{i=1}^m h_i(x).$

still reduces variance!

- idea: individual examples are still <sup>source</sup> distributed according to  $\mathcal{P}$  ~~dist~~
- even though they're not independent, they're "independent enough" to reduce variance.

Random forest  $\Rightarrow$  full algorithm Bagging + Trees



## General Gradient Boosted Regression Trees

Input: loss  $l$ ,  $\alpha$  step size,  $\mathcal{D}$ , CART alg

Init:  $H_0 = 0$ ;  $H(x) = 0$

Loop from  $t=1:T$

$r_i = l'(H_{t-1}(x_i); y_i)$  for all  $i \in \{1, \dots, n\}$

$h_t = \text{CART}(\{(x_1, -r_1), (x_2, -r_2), \dots, (x_n, -r_n)\})$

if  $\sum_{i=1}^n r_i h(x_i) < 0$ :

$H_T(x) = H_{T-1}(x) + \alpha h(x)$

else:

halt and return  $H_{t-1}$ .

end loop

return  $H_T$ .

could replace with  $l(H_{t-1} + \alpha h) < l(H_{t-1})$

## Random forest

Bagging + CART + Subsampling features

each node searches  
all features & splits

only search  $k$   
features at random



increases "diversity" of trees

set  $k \approx \sqrt{d}$

Hyper params:  $k \Leftarrow$  # of features we search

$m \Leftarrow$  # of datasets/trees  $\Rightarrow$  set as large as possible  
sampled with replacement

1. Sample  $m$  "bagged" datasets from  $\mathcal{D}$

call  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_m$

2. for each bagged dataset  $\mathcal{D}_i$

run a variant of CART/ID3

$\Rightarrow$  at each split choose  $k$  features at random and only consider splitting those  
outputting hypothesis  $h_i$

3. Output random forest classifier/regressor

$$\hat{h}(x) = \frac{1}{m} \sum_{i=1}^m h_i(x).$$



## Two Variants of Rf:

- build tree using one dataset and label the leaves with another independent dataset  
(prove consistency)
- build smaller trees (not to full depth)  
(pruning)

# sequential vs parallel Boosting vs Bassing

- sequentially call our tree-building algo, at each step building a tree that improves the ensemble!

Goal: create an ensemble classifier/regressor

$$H_T(x) = \sum_{t=1}^T \alpha_t h_t(x) = H_{T-1}(x) + \alpha_T h_T(x)$$

ensemble after T steps      scalar weights      hypothesis learned at step t

Goal: minimize loss function

$$l(H) = \frac{1}{n} \sum_{i=1}^n l(H(x_i), y_i)$$

the hypothesis I pick at step t+1 is

$$h_{t+1} = \arg \min_{h \in H} l(H_T + \alpha h)$$

~~$\alpha \in \mathbb{R}$~~  for simplicity imagine fixed  $\alpha$

$$= \arg \min_{h \in H} \frac{1}{n} \sum_{i=1}^n l(H(x_i) + \alpha h(x_i), y_i)$$

$\approx l(H(x_i), y_i) + \alpha \langle h(x_i), \nabla l(H(x_i), y_i) \rangle$

$$\approx l(H(x_i), y_i) + \alpha \langle h(x_i), \nabla l(H(x_i), y_i) \rangle$$

$$= \arg \min_{h \in H} \frac{1}{n} \sum_{i=1}^n \langle h(x_i), \nabla l(H(x_i), y_i) \rangle$$



$$= \operatorname{Arg\,min}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h(x_i) \ell'(H(x_i), y_i)$$

$$= \operatorname{Arg\,min}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h(x_i) \frac{\partial \ell(H(x_i), y_i)}{\partial [H(x_i)]}$$

$$\text{let } r_i = \frac{\partial \ell(H(x_i), y_i)}{\partial [H(x_i)]} = \ell'(H(x_i), y_i)$$

$$= \operatorname{Arg\,min}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h(x_i) \cdot r_i$$

$\nearrow \in \mathbb{R}$       $\nearrow \in \mathbb{R}$

if this is negative, then  $\ell(H_{T+1}) \leq \ell(H_T)$

Generic Boosting (a.k.a. "AnyBoost")

1. Assume that hypothesis  $h$  has constant

$$\sum_{i=1}^n (h(x_i))^2. \quad \text{e.g. if } h(x_i) \in \{+1, -1\}.$$

$$= \operatorname{Arg\,min}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (h(x_i) + r_i)^2 = \frac{1}{2} (h(x_i))^2 - \frac{1}{2} (r_i)^2$$

$$ab = \frac{1}{2} (a+b)^2 - \frac{1}{2} a^2 - \frac{1}{2} b^2$$

$$= \operatorname{Arg\,min}_{h \in \mathcal{H}} \frac{1}{2n} \sum_{i=1}^n (h(x_i) - (-r_i))^2$$

$t_i = -r_i = -\ell'(H(x_i), y_i)$

$$\approx \text{CART} \left( (x_1, t_1), (x_2, t_2), \dots, (x_n, t_n) \right)$$

# Ada Boost

Assumption: binary classification ( $y_i \in \{+1, -1\}$ )

→ weak learners also binary ( $h_i(x) \in \{+1, -1\}$ )

→ "step size"  $\alpha$  set by line search  
set optimally to minimize loss

→ loss Exponential loss  $\ell(H) = \sum_{i=1}^n \exp(-y_i H(x_i))$

Single ex loss:  $\ell(H(x); y) = \exp(-y H(x))$

$$\ell'(H(x); y) = -y \exp(-y H(x))$$

Boosting goal:  $\arg \min_{h \in \mathcal{H}} \sum_{i=1}^n \ell(H(x_i); y_i) \cdot h(x_i)$

$$= \arg \min_h - \sum_{i=1}^n y_i \exp(-y_i H(x_i)) \cdot h(x_i)$$

$$\text{let } w_i = \frac{1}{Z} \exp(-y_i H(x_i)), \quad Z = \sum_{i=1}^n \exp(-y_i H(x_i))$$

$$= \arg \min_h - \sum_{i=1}^n w_i y_i h(x_i)$$

$w_i \in \{-1, 1\}$   
binary  $\in \{-1, 1\}$

$$= \arg \min_h \sum_{i|y_i \neq h(x_i)} w_i - \sum_{i|y_i = h(x_i)} w_i$$

$1 - \sum_{i|y_i = h(x_i)} w_i$



$$= \arg \min_h 2 \sum_{i|y_i \neq h(x_i)} w_i - 1$$

$$= \arg \min_h \sum_{i|y_i \neq h(x_i)} w_i$$

weighted classification error of hypothesis  $h$

Choosing  $\alpha$ :

$$\text{pick } \alpha = \arg \min_{\alpha} \mathcal{L}(H + \alpha h)$$

$$= \arg \min_{\alpha} \sum_{i=1}^n \exp(-y_i (H(x_i) + \alpha h(x_i)))$$

$$= \arg \min_{\alpha} \sum_{i=1}^n \exp(-y_i H(x_i)) \exp(-y_i h(x_i) \alpha)$$

$$= \arg \min_{\alpha} \sum_{i=1}^n w_i \exp(-y_i h(x_i) \alpha)$$

$$= \arg \min_{\alpha} \sum_{i|y_i = h(x_i)} w_i \exp(-\alpha) + \sum_{i|y_i \neq h(x_i)} w_i \exp(+\alpha)$$

let  $\epsilon = \sum_{i|y_i \neq h(x_i)} w_i$   $\Leftarrow$  weighted class error of  $h$

$$= \arg \min_{\alpha} (1 - \epsilon) \exp(-\alpha) + \epsilon \exp(+\alpha)$$

$$= \alpha = \frac{1}{2} \log \left( \frac{1 - \epsilon}{\epsilon} \right)$$

Finishing up AdaBoost:

- incrementally update  $w_i$

$$w_i' = \frac{1}{Z'} \exp(-\gamma_i (H(x_i) + \alpha h(x_i)))$$

$$= \frac{Z}{Z'} w_i \cdot \exp(-\gamma_i h(x_i) \alpha) = \frac{w_i \exp(-\alpha \gamma_i h(x_i))}{2\sqrt{\epsilon(1-\epsilon)}}$$

$$Z' = \sum_{i=1}^n \exp(-\gamma_i (H(x_i) + \alpha h(x_i)))$$

$$= \sum_{i=1}^n Z \cdot w_i \cdot \exp(-\alpha \gamma_i h(x_i))$$

$$= Z \cdot \left( \sum_{\substack{i=1 \\ (|y_i \neq h(x_i))}} w_i \exp(\alpha) + \sum_{\substack{i=1 \\ (|y_i = h(x_i))}} w_i \exp(-\alpha) \right)$$

$$= Z \cdot (\epsilon \exp(\alpha) + (1-\epsilon) \exp(-\alpha))$$

recall:  $\alpha = \frac{1}{2} \log\left(\frac{1-\epsilon}{\epsilon}\right)$   $\exp(\alpha) = \sqrt{\frac{1-\epsilon}{\epsilon}}$

$$Z' = Z \cdot \left( \epsilon \cdot \sqrt{\frac{1-\epsilon}{\epsilon}} + (1-\epsilon) \cdot \sqrt{\frac{\epsilon}{1-\epsilon}} \right)$$

$$= Z \cdot (2\sqrt{\epsilon(1-\epsilon)})$$

$$\mathcal{L}(H+\alpha h) = \mathcal{L}(H) \cdot (2\sqrt{\epsilon(1-\epsilon)})$$





if we halt whenever the classifier  $h$   
has weighted test error  $\epsilon$  too  
close to  $\frac{1}{2}$ , i.e. suppose/require  
that  $\epsilon \leq \frac{1}{2} - \gamma$

$$\begin{aligned} \text{then: } 2\sqrt{\epsilon(1-\epsilon)} &\leq 2\sqrt{\left(\frac{1}{2}-\gamma\right)\left(\frac{1}{2}+\gamma\right)} \\ &= 2\left(\sqrt{\frac{1}{4}-\gamma^2}\right) \\ &= \sqrt{1-4\gamma^2} \end{aligned}$$

$$\mathcal{L}(H+\alpha h) \leq \mathcal{L}(H) \cdot \sqrt{1-4\gamma^2}$$

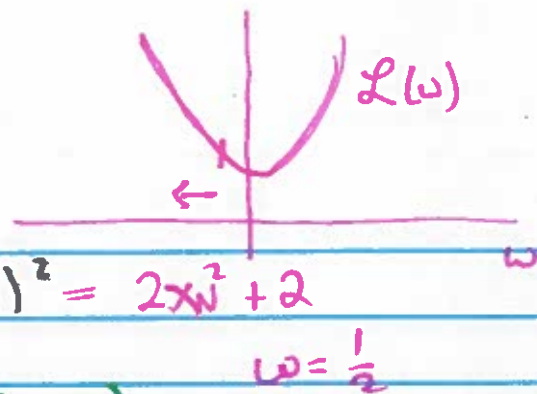
after we add  $k$  trees to the ensemble:

$$\mathcal{L}(H) \leq \mathcal{L}(O) \cdot (1-4\gamma^2)^{k/2} \leq n \cdot \underbrace{(1-4\gamma^2)^{k/2}}$$

if  $\gamma$  constant, need  $O(\log n)$  trees  
to achieve small loss

e.g.  $\mathcal{L}(w) = f_1(w) + f_2(w)$

$$= (w-1)^2 + (w+1)^2 = 2w^2 + 2$$



$$\mathcal{L}'(w) = 4w \quad f_1'(w) = 2(w-1)$$

$$\mathcal{L}'(w) f_1'(w) = 8w(w-1) \neq 0$$

Why SGD:  $(x_i, y_i) = \text{example set'd at time } t$

$$\mathcal{L}(w_{t+1}) = \mathcal{L}(w_t) - \alpha \nabla \mathcal{L}(w_t)^T \nabla \ell(h(x_i; w_t); y_i) + O(\alpha^2)$$

$$E[\mathcal{L}(w_{t+1})] = \mathcal{L}(w_t) - E[\alpha \nabla \mathcal{L}(w_t)^T \nabla \ell(h(x_i; w_t); y_i)] + O(\alpha^2)$$

$$= \mathcal{L}(w_t) - \alpha \nabla \mathcal{L}(w_t)^T E[\nabla \ell(h(x_i; w_t); y_i)] + O(\alpha^2)$$

$$E[\nabla \ell(h(x_i; w_t); y_i)] = \sum_{i=1}^n \left(\frac{1}{n}\right) \nabla \ell(h(x_i; w_t); y_i)$$

$$= \nabla_w \left( \frac{1}{n} \sum_{i=1}^n \ell(h(x_i; w_t); y_i) \right)$$

$\mathcal{L}(w_t)$

$$= \mathcal{L}(w_t) - \alpha \nabla \mathcal{L}(w_t)^T \nabla \mathcal{L}(w_t) + O(\alpha^2)$$

$$\frac{\|\nabla \mathcal{L}(w_t)\|^2}{2\alpha}$$

SGD vs GD:

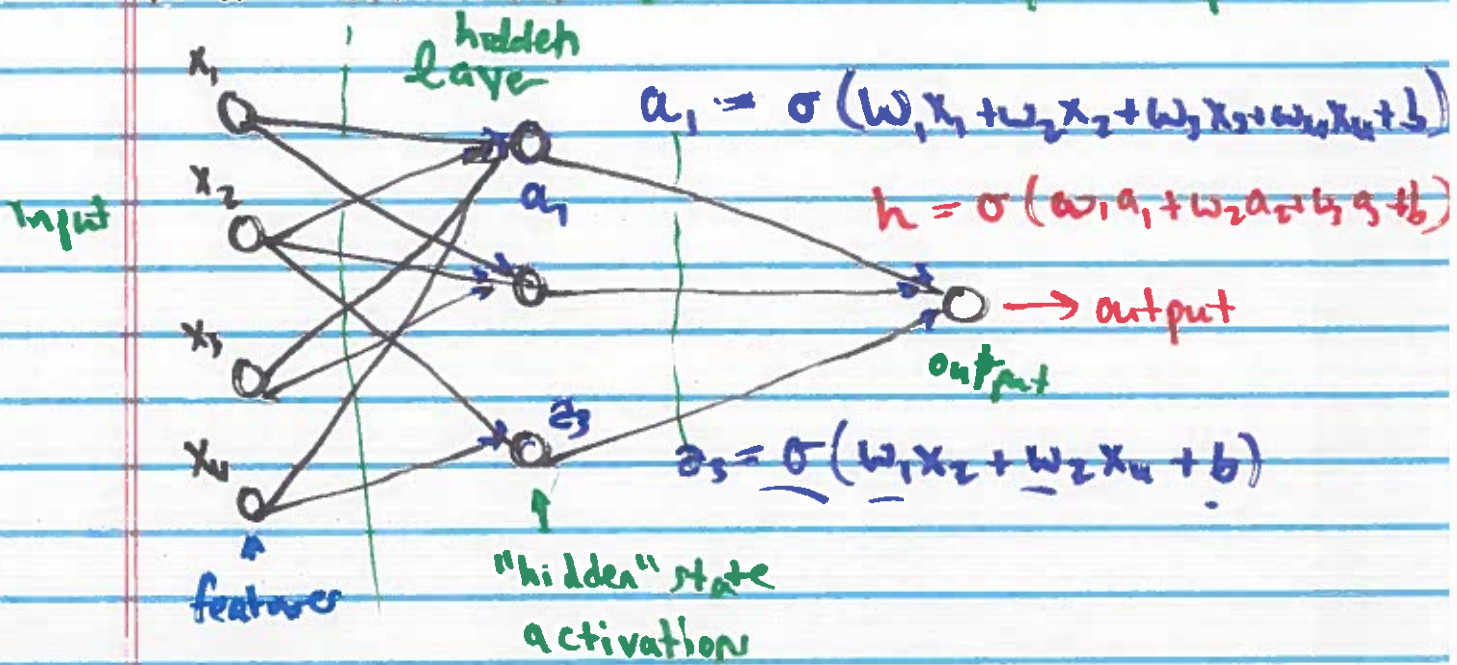
\* SGD is faster per iteration

\* SGD generalizes better!

\* SGD helps avoid local minima



# Neural Networks a.k.a "Multilayer Perceptron"



$\sigma$  acts elementwise

$$a \in \mathbb{R}^3$$

$$a = \sigma(Wx + b)$$

$\mathbb{R}^{3 \times 4} \quad \mathbb{R}^3$

$$\begin{aligned} a_1 &= \sigma(w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4 + b_1) \\ a_2 &= \sigma(\dots) \\ a_3 &= \sigma(\dots) \end{aligned}$$

$$h = \sigma(u_1a_1 + u_2a_2 + u_3a_3 + c)$$

$$h(x) = \sigma(u^T \sigma(Wx + b) + c)$$

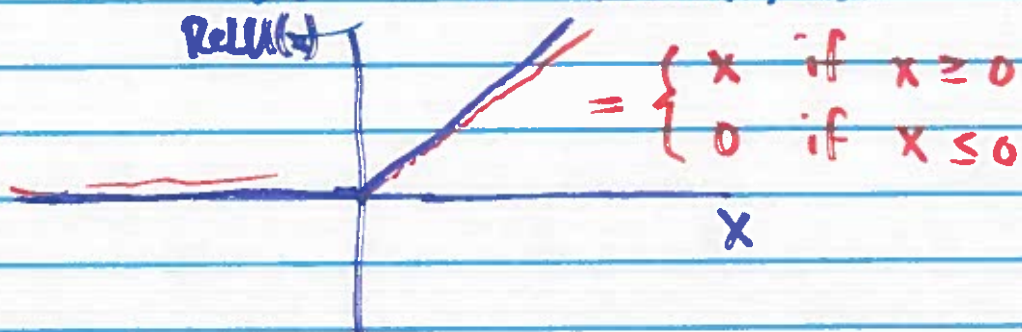
$\mathbb{R}^3 \quad \mathbb{R}^{3 \times 4} \quad \mathbb{R}^3 \quad \mathbb{R}$

$$h(x) = u^T \phi(x) + c, \text{ where } \phi(x) = \sigma(Wx + b) \in \mathbb{R}^3$$

What is  $\sigma$ ?

- ReLU: Rectified linear unit.

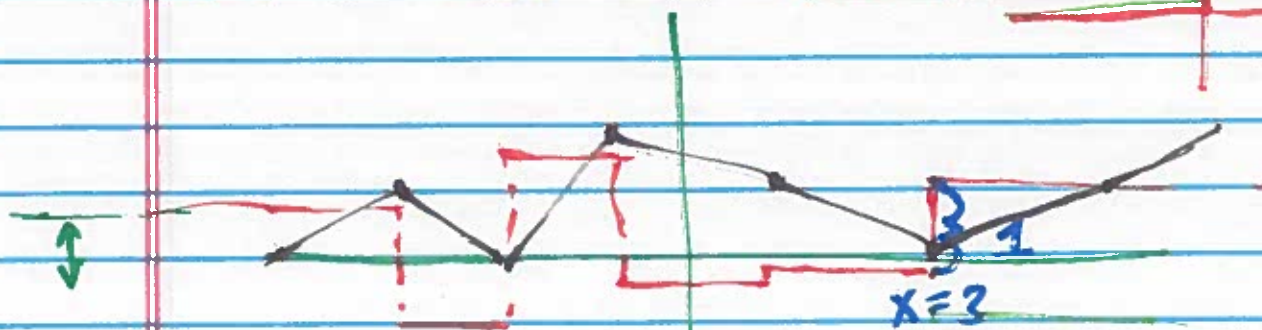
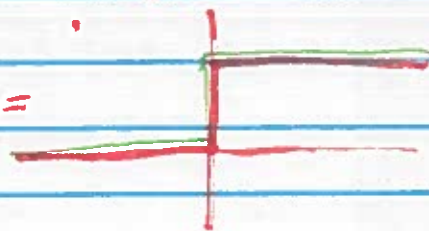
$$\text{ReLU}(x) = \sigma(x) = \max(0, x)$$



ReLU networks (nets where all nonlinearities are ReLU) are piecewise linear.

& all continuous piecewise linear <sup>on bounded domain</sup> functions can be expressed as a ReLU network.

Why does this hold?  $\text{ReLU}'(x) =$



$$1 \cdot \text{ReLU}'(x-3)$$

$$f'(x) = c_0 + a_1 \text{ReLU}'(x-b_1) + a_2 \text{ReLU}'(x-b_2) \dots$$

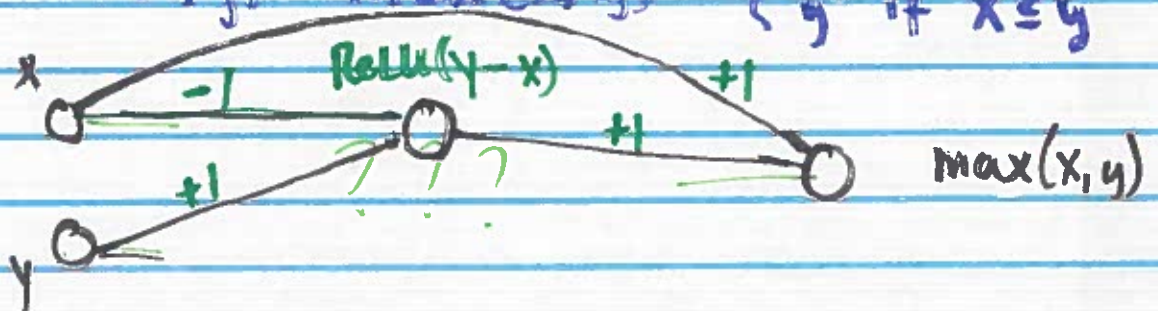
$$f(x) = c_0 x + c_1 + a_1 \text{ReLU}(x-b_1) + a_2 \text{ReLU}(x-b_2) \dots$$



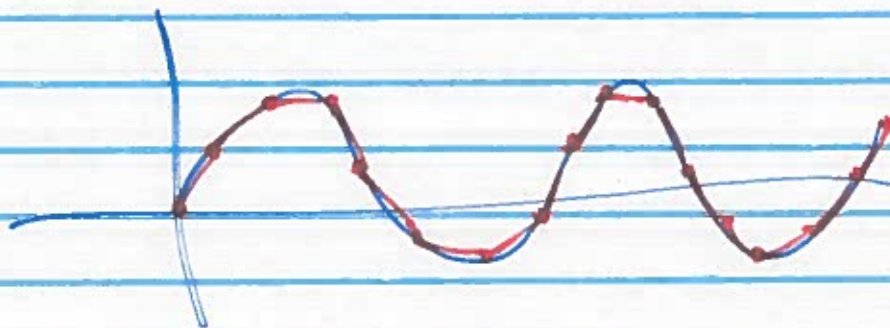
$$\text{Let } f(x) = |x| = \begin{cases} -x & \text{if } x \leq 0 \\ +x & \text{if } x \geq 0 \end{cases}$$

$$f(x) = -x + 2\text{ReLU}(x)$$

$$\text{Let } f(x, y) = \max(x, y) = \begin{cases} x & \text{if } x \geq y \\ y & \text{if } x \leq y \end{cases}$$



$$f(x, y) = \text{ReLU}(y-x) + x$$



any function  $\approx$  some piecewise linear function  $\equiv$  ANN + ReLU

"universality"  $\Rightarrow$  approx any function!

choose

learn

"architecture" + "weights/parameters"



$$w_1, w_2, \dots, b_1, \dots$$

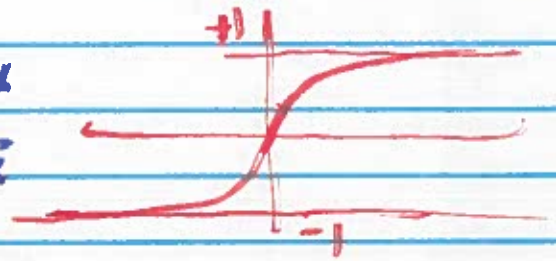
## Common Nonlinearities

- ReLU

- $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

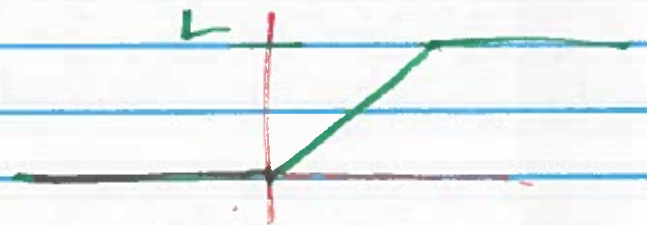
||

- $\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$



- "Smoothed" ReLU

- "saturation" ReLU



$$\sigma(x) = \max(0, \min(x, L))$$



How to learn a DNN:

STOCHASTIC  
Gradient descent! ) Using a subsample  
of  $\mathcal{D}$  to  
compute  $\nabla$  gradients  
at each step

$$\min_w \sum_{i=1}^n \ell(h(x_i; w); y_i)$$

+ automatic differentiation

\* not numerical differentiation  
 $f'(x) \approx \frac{f(x+0.01) - f(x)}{0.01}$

\* not symbolic diff

\* about the same time as computing  
 $f$  to compute  $\nabla f$

→ BACKPROPAGATION ↗

$$\text{GD: } w \leftarrow w - \alpha \frac{1}{n} \sum_{i=1}^n \nabla \ell(h(x_i; w), y_i)$$

SGD: pick  $i$  from  $\{1, \dots, n\}$  unif at random  
 $w \leftarrow w - \alpha \nabla \ell(h(x_i; w), y_i)$

$$\mathbb{E} \ell(\underbrace{w - \alpha \nabla \ell(w)}_{\text{next } w}) = \ell(w) - \alpha \nabla \ell(w)^T \nabla \ell(w) + O(\alpha^2)$$

$\|\nabla \ell(w)\|^2$   
 $\geq 0$

$$\mathbb{E} \ell(w - \nabla \ell(h(x_i; w), y_i)) = \ell(w) - \alpha \nabla \ell(w)^T \nabla \ell(h(x_i; w), y_i) + O(\alpha^2)$$

# Convolutional Neural Network.

Image is what? Multidim array ( $3 \times W \times H$ )

e.g. a  $16 \times 16$  color img: ( $3 \times 16 \times 16$ ) array

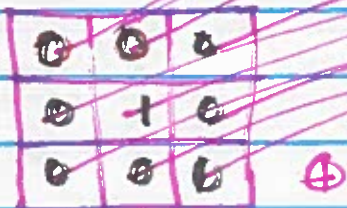
$$h(x) = U \cdot \text{ReLU}(Wx + b)$$

$$= U \cdot \text{ReLU}(W \cdot \text{vec}(x) + b)$$

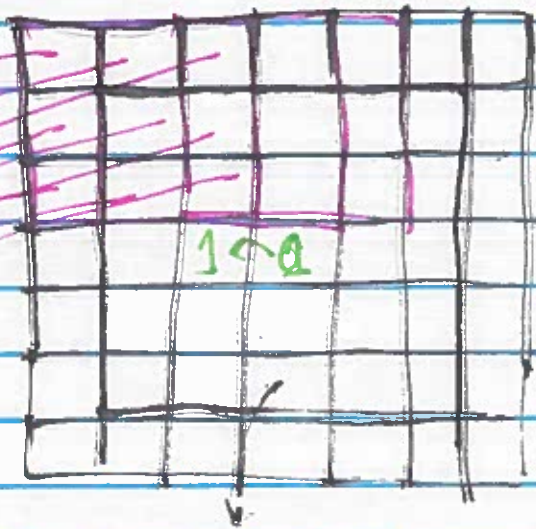
$\mathbb{R}^{3 \cdot 16 \cdot 16}$

Want: something shift-invariant

convolution!



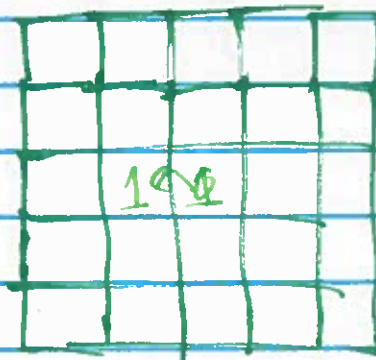
filter



1 ⊕

$$\text{out} = \text{conv}(\text{in}, \text{filter})$$

$$(d-k) \times (d-k) \quad d \times d \quad k \times k$$

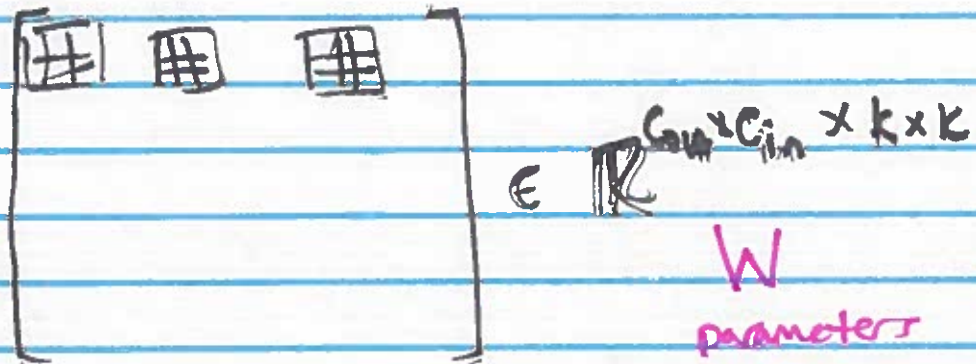


1 ⊕

5x5

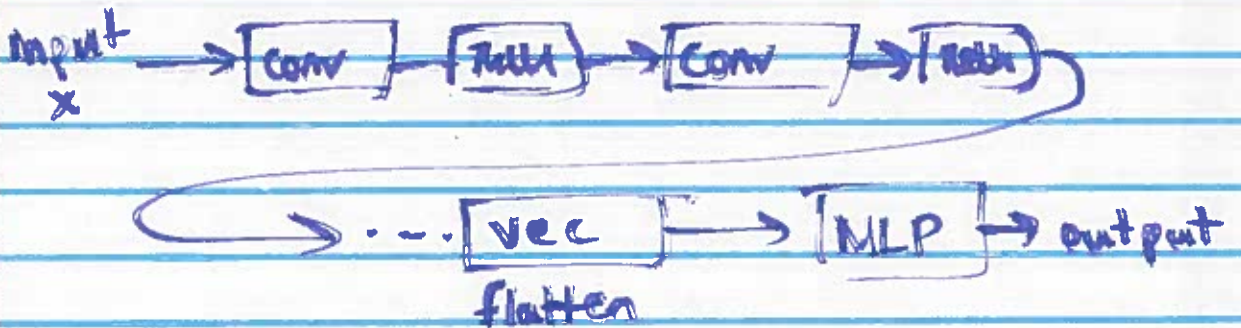


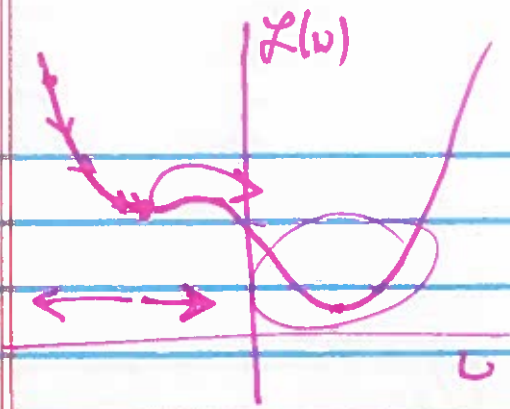
Conv layers:



$$(out)_j = \sum_i conv((in)_i, W_{j,i})$$

CNN





## DNN and Overfitting

a model is overparameterized when  $d \gg n$

number of parameters

number of training examples



"DNNs don't overfit" \*

as much as other models when trained using standard methods

Counter-overfitting methods:

\*  $l_2$  regularization  $+ \lambda \|w\|^2$  "weight decay"

\* dropout: while training, for each step of SGD, "remove" each hidden node from the network w.p.  $p > 0$ , (e.g.  $p = \frac{1}{2}$ )

\* batch normalization