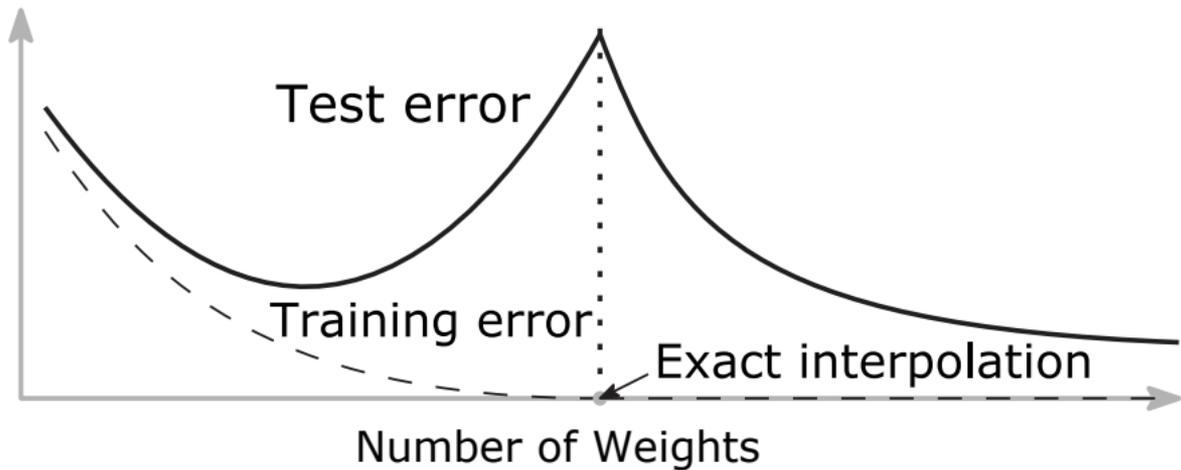


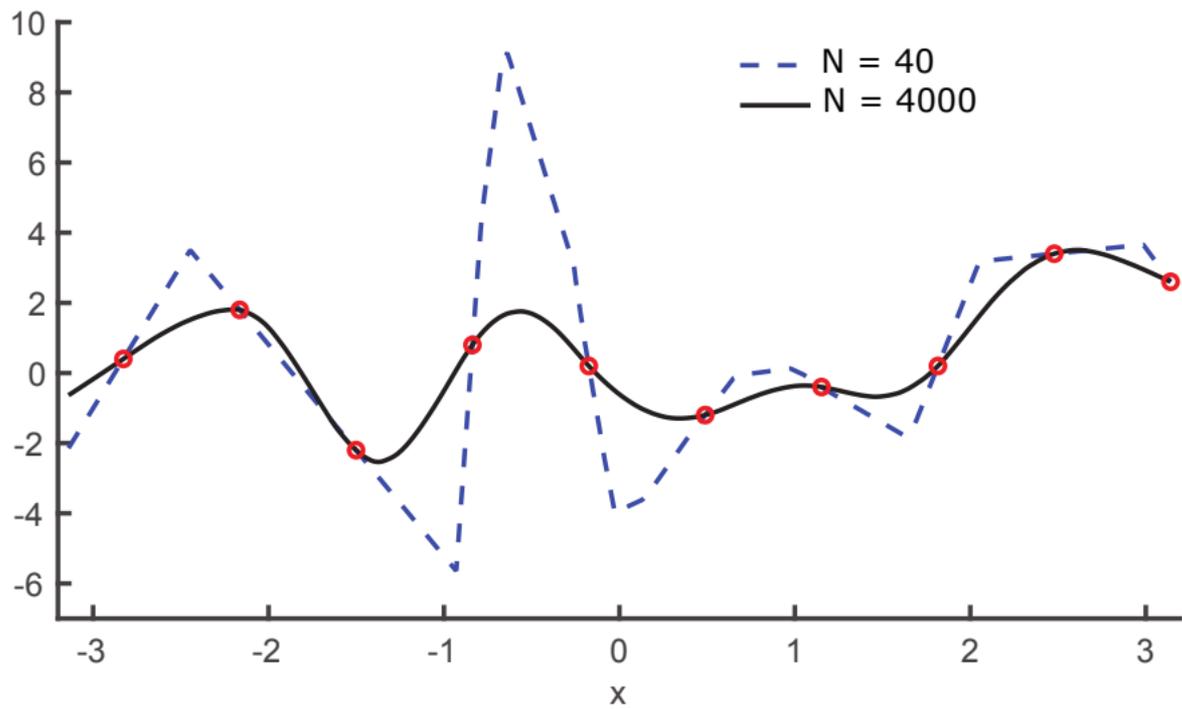
# Deep Learning/Double Descent

**Gilbert Strang**

**MIT**

**October, 2019**





# Fit training data by a Learning function $F$

We are given training data: **Inputs**  $v$ , **outputs**  $w$

**Example** Each  $v$  is an image of a number  
 $w = 0, 1, \dots, 9$

The vector  $v$  describes each pixel in the image

We want to create a **learning function** so that

$$F(v) \approx w$$

[math.mit.edu/learningfromdata](http://math.mit.edu/learningfromdata)

# Deep Neural Networks

- 1 Key operation      Composition  $F = F_3(F_2(F_1(x, v_0)))$
- 2 Key rule            Chain rule for  $x$ -derivatives of  $F$
- 3 Key algorithm      Stochastic gradient descent to find  $x$
- 4 Key subroutine     Backpropagation to compute grad  $F$
- 5 Key nonlinearity    ReLU  $(y) = \max(y, 0) = \text{ramp function}$

# Deep Neural Networks

- 1 Key operation      Composition  $F = F_3(F_2(F_1(x, v_0)))$
- 2 Key rule            Chain rule for  $x$ -derivatives of  $F$
- 3 Key algorithm      Stochastic gradient descent to find  $x$
- 4 Key subroutine     Backpropagation to compute grad  $F$
- 5 Key nonlinearity    ReLU  $(y) = \max(y, 0) = \text{ramp function}$

**Layer  $k$**      $v_k = F_k(v_{k-1}) = \text{ReLU}(A_k v_{k-1} + b_k)$

# Deep Neural Networks

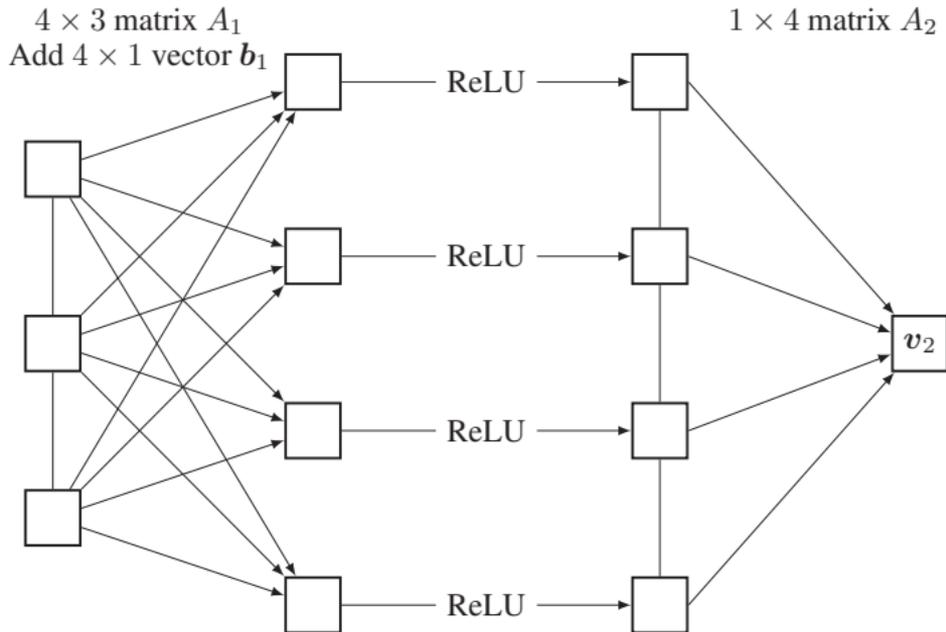
- 1 Key operation      Composition  $F = F_3(F_2(F_1(x, v_0)))$
- 2 Key rule            Chain rule for  $x$ -derivatives of  $F$
- 3 Key algorithm      Stochastic gradient descent to find  $x$
- 4 Key subroutine    Backpropagation to compute grad  $F$
- 5 Key nonlinearity   ReLU  $(y) = \max(y, 0) = \text{ramp function}$

**Layer  $k$**      $v_k = F_k(v_{k-1}) = \text{ReLU}(A_k v_{k-1} + b_k)$

Weights for layer  $k$      $A_k = \text{matrix}$  and  $b_k = \text{offset vector}$

$v_0 = \text{training data}$  /  $v_1, \dots, v_{\ell-1}$  hidden layers /  $v_\ell = \text{output}$

# Figure from [math.mit.edu/learningfromdata](http://math.mit.edu/learningfromdata)



Feature vector  $\mathbf{v}_0$   
Three components for  
each training sample

$\mathbf{y}_1 = A_1 \mathbf{v}_0 + \mathbf{b}_1$   
 $\mathbf{y}_1$  at layer 1

Four components of  $\mathbf{y}_1$  and  $\mathbf{v}_1$

$\mathbf{v}_1$  at layer 1  
 $\mathbf{v}_1 = \text{ReLU}(\mathbf{y}_1)$

Output  $\mathbf{w} = \mathbf{v}_2$   
 $\mathbf{v}_2 = A_2 \mathbf{v}_1$

**Key computation:** Weights  $\boldsymbol{x}$  minimize overall loss  $L(\boldsymbol{x})$

---

$$L(\boldsymbol{x}) = \frac{1}{N} \sum_{j=1}^N \text{loss } \ell(\boldsymbol{x}, \boldsymbol{v}_0^j) \text{ on sample } j$$

## Key computation: Weights $x$ minimize overall loss $L(x)$

$$L(x) = \frac{1}{N} \sum_{j=1}^N \text{loss } \ell(x, v_0^j) \text{ on sample } j$$

“Square loss” = error  $\ell(x, v_0^j) = \left\| F(x, v_0^j) - \text{true} \right\|^2$

Cross-entropy loss, hinge loss, ...

Classification problem: true = 1 or -1

Regression problem: true = vector

Gradient descent  $x_{k+1} = \arg \min \|x_k - s_k \nabla L(x_k)\|$

Stochastic descent  $x_{k+1} = \arg \min \|x_k - s_k \nabla \ell(x_k, v_0^k)\|$

# Key Questions

---

1. Optimization of the weights  $\mathbf{x} = A_k$  and  $\mathbf{b}_k$

# Key Questions

---

1. Optimization of the weights  $\mathbf{x} = A_k$  and  $\mathbf{b}_k$
2. Convergence rate of descent and accelerated descent  
(when  $\mathbf{x}_{k+1}$  depends on  $\mathbf{x}_k$  and  $\mathbf{x}_{k-1}$ : momentum added)

# Key Questions

---

1. Optimization of the weights  $\boldsymbol{x} = A_k$  and  $\boldsymbol{b}_k$
2. Convergence rate of descent and accelerated descent  
(when  $\boldsymbol{x}_{k+1}$  depends on  $\boldsymbol{x}_k$  and  $\boldsymbol{x}_{k-1}$ : momentum added)
3. Do the weights  $A_1, b_1 \dots$  generalize to unseen test data?  
(Early stopping / Do not overfit the data)

- 
1. Stochastic gradient descent optimizes weights  $A_k, b_k$

1. Stochastic gradient descent optimizes weights  $A_k, b_k$
2. Backpropagation in the computational graph computes derivatives with respect to weights  $x = A_1, b_1, \dots, A_\ell, b_\ell$

1. Stochastic gradient descent optimizes weights  $A_k, b_k$
2. Backpropagation in the computational graph computes derivatives with respect to weights  $x = A_1, b_1, \dots, A_\ell, b_\ell$
3. The learning function  $F(\mathbf{x}, \mathbf{v}_0) = \dots F_3(F_2(F_1(\mathbf{x}, \mathbf{v})))$

1. Stochastic gradient descent optimizes weights  $A_k, b_k$
2. Backpropagation in the computational graph computes derivatives with respect to weights  $x = A_1, b_1, \dots, A_\ell, b_\ell$
3. The learning function  $F(\mathbf{x}, \mathbf{v}_0) = \dots F_3(F_2(F_1(\mathbf{x}, \mathbf{v})))$

$$F_1(\mathbf{v}_0) = \max(A_1 \mathbf{v}_0 + b_1, 0) = \text{ReLU} \circ \text{affine map}$$

$F(\mathbf{v})$  is continuous piecewise linear : how many pieces?

This measures the “expressivity” of the network

Assume 1 hidden layer with  $N$  neurons

---

$v_0$  has  $m$  components /  $v_1$  has  $N$  components /  $N$  ReLU's

---

$v_0$  has  $m$  components /  $v_1$  has  $N$  components /  $N$  ReLU's

The number of flat regions in  $\mathbf{R}^m$  bounded by the  $N$  hyperplanes

$$r(N, m) = \sum_{i=0}^m \binom{N}{i} = \binom{N}{0} + \binom{N}{1} + \cdots + \binom{N}{m}$$

---

$v_0$  has  $m$  components /  $v_1$  has  $N$  components /  $N$  ReLU's

The number of flat regions in  $\mathbf{R}^m$  bounded by the  $N$  hyperplanes

$$r(N, m) = \sum_{i=0}^m \binom{N}{i} = \binom{N}{0} + \binom{N}{1} + \dots + \binom{N}{m}$$

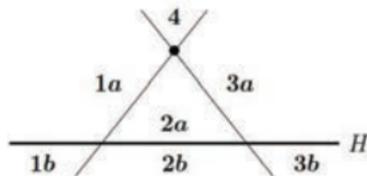
$N = 3$  folds in a plane will produce  $1 + 3 + 3 = 7$  pieces

$v_0$  has  $m$  components /  $v_1$  has  $N$  components /  $N$  ReLU's

The number of flat regions in  $\mathbf{R}^m$  bounded by the  $N$  hyperplanes

$$r(N, m) = \sum_{i=0}^m \binom{N}{i} = \binom{N}{0} + \binom{N}{1} + \dots + \binom{N}{m}$$

$N = 3$  folds in a plane will produce  $1 + 3 + 3 = 7$  pieces



Start with 2 folds

$$\leftarrow r(2, 2) = 4$$

Add new fold

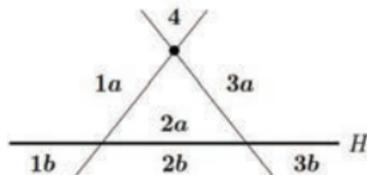
$$\leftarrow r(2, 1) = 3$$

$v_0$  has  $m$  components /  $v_1$  has  $N$  components /  $N$  ReLU's

The number of flat regions in  $\mathbf{R}^m$  bounded by the  $N$  hyperplanes

$$r(N, m) = \sum_{i=0}^m \binom{N}{i} = \binom{N}{0} + \binom{N}{1} + \dots + \binom{N}{m}$$

$N = 3$  folds in a plane will produce  $1 + 3 + 3 = 7$  pieces



Start with 2 folds

$$\leftarrow r(2, 2) = 4$$

Add new fold

$$\leftarrow r(2, 1) = 3$$

$$\text{Recursion } r(N, m) = r(N - 1, m) + r(N - 1, m - 1)$$

## **$F(x) = F_2(F_1(x))$ is continuous piecewise linear**

---

One hidden layer of neurons: deep networks have many more

## **$F(x) = F_2(F_1(x))$ is continuous piecewise linear**

---

One hidden layer of neurons: deep networks have many more

Overfitting is not desirable! Gradient descent stops early!

## **$F(x) = F_2(F_1(x))$ is continuous piecewise linear**

---

One hidden layer of neurons: deep networks have many more

Overfitting is not desirable! Gradient descent stops early!

“Generalization” measured by success on unseen test data

## $F(x) = F_2(F_1(x))$ is continuous piecewise linear

One hidden layer of neurons: deep networks have many more

Overfitting is not desirable! Gradient descent stops early!

“Generalization” measured by success on unseen test data

Big problems are **underdetermined** [# weights > # samples]

## $F(x) = F_2(F_1(x))$ is continuous piecewise linear

One hidden layer of neurons: deep networks have many more

Overfitting is not desirable! Gradient descent stops early!

“Generalization” measured by success on unseen test data

Big problems are **underdetermined** [# weights > # samples]

Stochastic Gradient Descent finds weights that generalize well

# Convolutional Neural Nets (**CNN**)

$$A = \begin{bmatrix} x_1 & x_0 & x_{-1} & 0 & 0 & 0 \\ 0 & x_1 & x_0 & x_{-1} & 0 & 0 \\ 0 & 0 & x_1 & x_0 & x_{-1} & 0 \\ 0 & 0 & 0 & x_1 & x_0 & x_{-1} \end{bmatrix}$$

$N + 2$  inputs and  $N$  outputs

# Convolutional Neural Nets (**CNN**)

$$A = \begin{bmatrix} x_1 & x_0 & x_{-1} & 0 & 0 & 0 \\ 0 & x_1 & x_0 & x_{-1} & 0 & 0 \\ 0 & 0 & x_1 & x_0 & x_{-1} & 0 \\ 0 & 0 & 0 & x_1 & x_0 & x_{-1} \end{bmatrix}$$

$N + 2$  inputs and  $N$  outputs

**Each shift has a diagonal of 1's**

$$A = x_1L + x_0C + x_{-1}R$$

# Convolutional Neural Nets (**CNN**)

$$A = \begin{bmatrix} x_1 & x_0 & x_{-1} & 0 & 0 & 0 \\ 0 & x_1 & x_0 & x_{-1} & 0 & 0 \\ 0 & 0 & x_1 & x_0 & x_{-1} & 0 \\ 0 & 0 & 0 & x_1 & x_0 & x_{-1} \end{bmatrix}$$

$N + 2$  inputs and  $N$  outputs

Each shift has a diagonal of 1's

$$A = x_1 L + x_0 C + x_{-1} R$$

$$\frac{\partial \mathbf{y}}{\partial x_1} = L \mathbf{v} \quad \frac{\partial \mathbf{y}}{\partial x_0} = C \mathbf{v} \quad \frac{\partial \mathbf{y}}{\partial x_{-1}} = R \mathbf{v}$$

# Convolutions in Two Dimensions

**Weights**  $\begin{bmatrix} x_{11} & x_{01} & x_{-11} \\ x_{10} & x_{00} & x_{-10} \\ x_{1-1} & x_{0-1} & x_{-1-1} \end{bmatrix}$  **Input image**  $v_{ij}$   $i, j$  from  $(0, 0)$  to  $(N+1, N+1)$   
**Output image**  $y_{ij}$   $i, j$  from  $(1, 1)$  to  $(N, N)$   
**Shifts L, C, R, U, D = Left, Center, Right, Up, Down**

# Convolutions in Two Dimensions

$$\text{Weights} \begin{bmatrix} x_{11} & x_{01} & x_{-11} \\ x_{10} & x_{00} & x_{-10} \\ x_{1-1} & x_{0-1} & x_{-1-1} \end{bmatrix} \begin{array}{l} \text{Input image } v_{ij} \text{ } i, j \text{ from } (0, 0) \text{ to } (N+1, N+1) \\ \text{Output image } y_{ij} \text{ } i, j \text{ from } (1, 1) \text{ to } (N, N) \\ \text{Shifts } \mathbf{L, C, R, U, D} = \mathbf{Left, Center, Right, Up, Down} \end{array}$$

A convolution is a combination of shift matrices = filter = Toeplitz matrix

The coefficients in the combination will be the “weights” to be learned.

# Convolutions in Two Dimensions

$$\text{Weights} \begin{bmatrix} x_{11} & x_{01} & x_{-11} \\ x_{10} & x_{00} & x_{-10} \\ x_{1-1} & x_{0-1} & x_{-1-1} \end{bmatrix} \begin{array}{l} \text{Input image } v_{ij} \text{ } i, j \text{ from } (0, 0) \text{ to } (N+1, N+1) \\ \text{Output image } y_{ij} \text{ } i, j \text{ from } (1, 1) \text{ to } (N, N) \\ \text{Shifts } \mathbf{L, C, R, U, D} = \mathbf{Left, Center, Right, Up, Down} \end{array}$$

A convolution is a combination of shift matrices = filter = Toeplitz matrix

The coefficients in the combination will be the “weights” to be learned.

**9 weights instead of  $n^2$  weights in  $A$**

**Allows the neural net to have more width and depth**

Computing the weights  $\boldsymbol{x}$  = matrices  $A_k$ , bias vectors  $b_k$

---

**Choose a loss function  $\ell$  to measure  $F(\boldsymbol{x}, v)$  – true output**

## Computing the weights $\boldsymbol{x}$ = matrices $A_k$ , bias vectors $b_k$

Choose a loss function  $\ell$  to measure  $F(\boldsymbol{x}, \boldsymbol{v})$  – true output

**Total loss**  $L = \frac{1}{N}$  (sum of losses for all  $N$  samples)

# Computing the weights $\boldsymbol{x}$ = matrices $A_k$ , bias vectors $b_k$

Choose a loss function  $\ell$  to measure  $F(\boldsymbol{x}, \boldsymbol{v})$  – true output

Total loss  $L = \frac{1}{N}$  (sum of losses for all  $N$  samples)

Compute weights  $\boldsymbol{x}$  to minimize the total loss  $L$

---

Here are three loss functions—Cross-entropy is a favorite loss function for neural nets

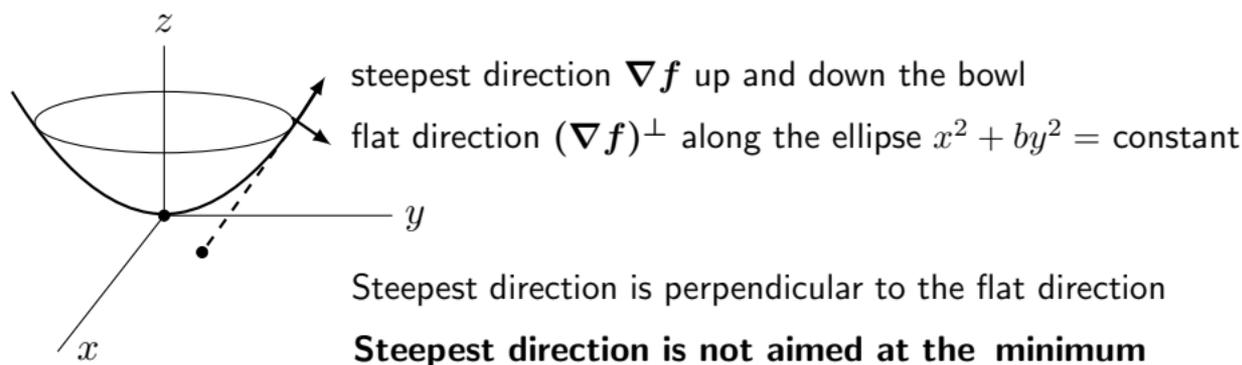
Here are three loss functions—Cross-entropy is a favorite loss function for neural nets

1 **Square loss**  $L(\mathbf{x}) = \frac{1}{N} \sum_1^N \|F(\mathbf{x}, \mathbf{v}_i) - \text{true}\|^2$  : sum over samples  $\mathbf{v}_i$

2 **Hinge loss**  $L(\mathbf{x}) = \frac{1}{N} \sum_1^N \max(0, 1 - t F(\mathbf{x}))$  for **classification**

3 **Cross-entropy loss**  $L(\mathbf{x}) = -\frac{1}{N} \sum_1^N [y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)]$

# Steepest Descent = Gradient Descent $f = x^2 + by^2$

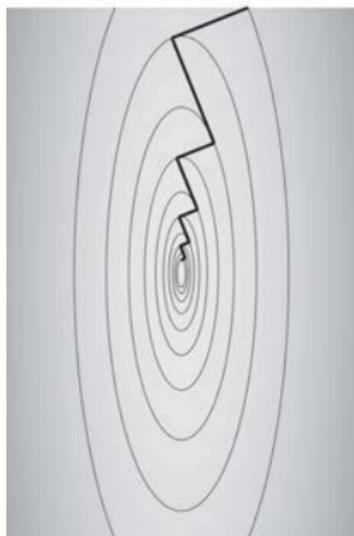


**Steepest descent moves in the gradient direction**  $\begin{bmatrix} -2x \\ -2by \end{bmatrix}$ .

$$x_k = b \left( \frac{b-1}{b+1} \right)^k \quad y_k = \left( \frac{1-b}{1+b} \right)^k \quad f(x_k, y_k) = \left( \frac{1-b}{1+b} \right)^{2k} f(x_0, y_0)$$

Descent formula  $\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \nabla F(\mathbf{x})$     Step size  $s_k =$  Learning rate

### Gradient Descent



The first descent step starts out perpendicular to the level set. As it crosses through lower level sets, the function  $f(x, y)$  is decreasing. **Eventually its path is tangent to a level set  $L$ .**

Slow convergence on a zig-zag path to the minimum of  $f = x^2 + by^2$ .

# Momentum and the Path of a Heavy Ball

**Descent with  
momentum**

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - s\mathbf{z}_k \\ \text{with } \mathbf{z}_k &= \nabla f(\mathbf{x}_k) + \beta\mathbf{z}_{k-1} \end{aligned}$$

**Descent with  
momentum**

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - s\mathbf{z}_k \\ \mathbf{z}_{k+1} - \nabla f(\mathbf{x}_{k+1}) &= \beta\mathbf{z}_k \end{aligned}$$

# Momentum and the Path of a Heavy Ball

Descent with  
momentum

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - s z_k \\ \text{with } z_k &= \nabla f(\mathbf{x}_k) + \beta z_{k-1} \end{aligned}$$

Descent with  
momentum

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - s z_k \\ z_{k+1} - \nabla f(\mathbf{x}_{k+1}) &= \beta z_k \end{aligned}$$

Following the  
eigenvector  $q$

$$\begin{aligned} c_{k+1} &= c_k - s d_k \\ -\lambda c_{k+1} + d_{k+1} &= \beta d_k \end{aligned} \quad \begin{bmatrix} \mathbf{1} & 0 \\ -\lambda & \mathbf{1} \end{bmatrix} \begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix}$$

# Momentum and the Path of a Heavy Ball

Descent with  
momentum

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - s \mathbf{z}_k \\ \text{with } \mathbf{z}_k &= \nabla f(\mathbf{x}_k) + \beta \mathbf{z}_{k-1} \end{aligned}$$

Descent with  
momentum

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - s \mathbf{z}_k \\ \mathbf{z}_{k+1} - \nabla f(\mathbf{x}_{k+1}) &= \beta \mathbf{z}_k \end{aligned}$$

Following the  
eigenvector  $\mathbf{q}$

$$\begin{aligned} c_{k+1} &= c_k - s d_k \\ -\lambda c_{k+1} + d_{k+1} &= \beta d_k \end{aligned} \quad \begin{bmatrix} \mathbf{1} & 0 \\ -\lambda & \mathbf{1} \end{bmatrix} \begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix}$$

It seems a miracle that this problem has a beautiful solution.

The optimal  $s$  and  $\beta$  are

$$s = \left( \frac{2}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}} \right)^2 \quad \text{and} \quad \beta = \left( \frac{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}} \right)^2$$

Key difference:  $b$  is replaced by  $\sqrt{b}$

**Ordinary  
descent factor**  $\left(\frac{1-b}{1+b}\right)^2$

**Accelerated  
descent factor**  $\left(\frac{1-\sqrt{b}}{1+\sqrt{b}}\right)^2$

# Key difference: $b$ is replaced by $\sqrt{b}$

Ordinary  
descent factor  $\left(\frac{1-b}{1+b}\right)^2$

Accelerated  
descent factor  $\left(\frac{1-\sqrt{b}}{1+\sqrt{b}}\right)^2$

**Steepest  
descent**  $\left(\frac{.99}{1.01}\right)^2 = .96$

**Accelerated  
descent**  $\left(\frac{.9}{1.1}\right)^2 = .67$

# Key difference: $b$ is replaced by $\sqrt{b}$

<b>Ordinary descent factor</b>	$\left(\frac{1-b}{1+b}\right)^2$
------------------------------------	----------------------------------

<b>Accelerated descent factor</b>	$\left(\frac{1-\sqrt{b}}{1+\sqrt{b}}\right)^2$
---------------------------------------	--

<b>Steepest descent</b>	$\left(\frac{.99}{1.01}\right)^2 = .96$
-----------------------------	---

<b>Accelerated descent</b>	$\left(\frac{.9}{1.1}\right)^2 = .67$
--------------------------------	---------------------------------------

Notice that  $\lambda_{\max}/\lambda_{\min} = 1/b = \kappa$  is the **condition number of  $S$**

# Stochastic Gradient Descent

---

**Stochastic gradient descent uses a “minibatch” of the training data**

# Stochastic Gradient Descent

---

Stochastic gradient descent uses a “minibatch” of the training data

Every step is much faster than using all data

# Stochastic Gradient Descent

---

Stochastic gradient descent uses a “minibatch” of the training data

Every step is much faster than using all data

*We don't want to fit the training data too perfectly (**overfitting**)*

# Stochastic Gradient Descent

---

Stochastic gradient descent uses a “minibatch” of the training data

Every step is much faster than using all data

*We don't want to fit the training data too perfectly (overfitting)*

*Choosing a polynomial of degree 60 to fit 61 data points*

# Stochastic Descent Using One Sample Per Step

---

Early steps of SGD often converge quickly toward the solution  $x^*$

# Stochastic Descent Using One Sample Per Step

Early steps of SGD often converge quickly toward the solution  $x^*$

Here we pause to look at semi-convergence: Fast start by stochastic gradient descent

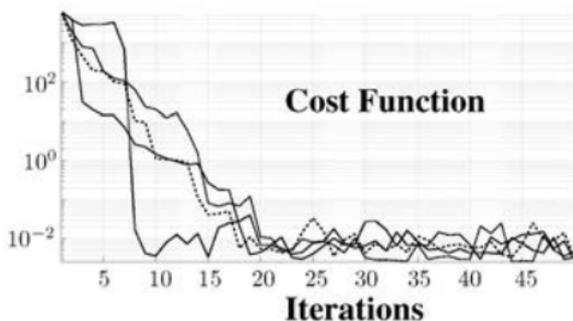
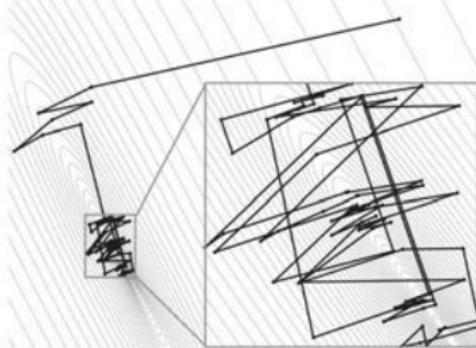
**Convergence at the start changes to large oscillations near the solution**

# Stochastic Descent Using One Sample Per Step

Early steps of SGD often converge quickly toward the solution  $x^*$

Here we pause to look at semi-convergence: Fast start by stochastic gradient descent

Convergence at the start changes to large oscillations near the solution

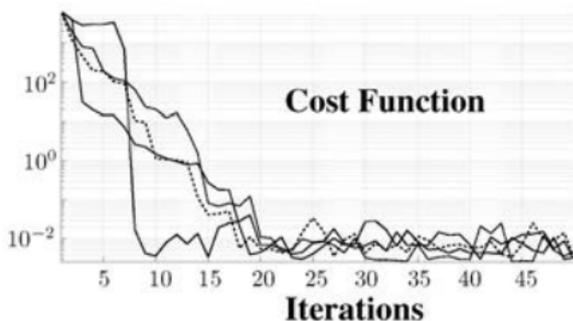
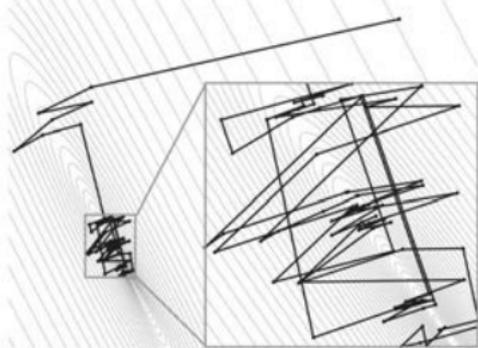


# Stochastic Descent Using One Sample Per Step

Early steps of SGD often converge quickly toward the solution  $x^*$

Here we pause to look at semi-convergence: Fast start by stochastic gradient descent

Convergence at the start changes to large oscillations near the solution



Kaczmarz for  $Ax = b$  with random  $i(k)$  
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{b_i - \mathbf{a}_i^T \mathbf{x}_k}{\|\mathbf{a}_i\|^2} \mathbf{a}_i$$

**Adaptive Stochastic Gradient Descent**

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \mathbf{D}_k$$

# Adaptive Methods Using Earlier Gradients (**ADAM**)

Adaptive Stochastic Gradient Descent

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \mathbf{D}_k$$

$$\mathbf{D}_k = \delta \mathbf{D}_{k-1} + (1 - \delta) \nabla L(\mathbf{x}_k) \quad s_k^2 = \beta s_{k-1}^2 + (1 - \beta) \|\nabla L(\mathbf{x}_k)\|^2$$

# Adaptive Methods Using Earlier Gradients (**ADAM**)

**Adaptive Stochastic Gradient Descent**

$$\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \mathbf{D}_k$$

$$\mathbf{D}_k = \delta \mathbf{D}_{k-1} + (1 - \delta) \nabla L(\mathbf{x}_k) \quad s_k^2 = \beta s_{k-1}^2 + (1 - \beta) \|\nabla L(\mathbf{x}_k)\|^2$$

**Why do the weights generalize well to unseen test data?**

# Computation of $\partial F / \partial x$ : Explicit Formulas

---

$$\mathbf{v}_L = \mathbf{b}_L + A_L \mathbf{v}_{L-1} \quad \text{or simply} \quad \mathbf{w} = \mathbf{b} + A\mathbf{v}.$$

# Computation of $\partial F / \partial x$ : Explicit Formulas

$$v_L = b_L + A_L v_{L-1} \quad \text{or simply} \quad w = b + Av.$$

*The output  $w_i$  is not affected by  $b_j$  or  $A_{jk}$  if  $j \neq i$*

# Computation of $\partial F / \partial x$ : Explicit Formulas

$$v_L = b_L + A_L v_{L-1} \quad \text{or simply} \quad w = b + Av.$$

*The output  $w_i$  is not affected by  $b_j$  or  $A_{jk}$  if  $j \neq i$*

**Fully connected layer**

**Independent weights  $A_{jk}$**

$$\frac{\partial w_i}{\partial b_j} = \delta_{ij} \quad \text{and} \quad \frac{\partial w_i}{\partial A_{jk}} = \delta_{ij} v_k$$

# Computation of $\partial F / \partial x$ : Explicit Formulas

$$v_L = b_L + A_L v_{L-1} \quad \text{or simply} \quad w = b + Av.$$

*The output  $w_i$  is not affected by  $b_j$  or  $A_{jk}$  if  $j \neq i$*

**Fully connected layer**

**Independent weights  $A_{jk}$**

$$\frac{\partial w_i}{\partial b_j} = \delta_{ij} \quad \text{and} \quad \frac{\partial w_i}{\partial A_{jk}} = \delta_{ij} v_k$$

Example 
$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} + \begin{bmatrix} a_{11}v_1 + a_{12}v_2 \\ a_{21}v_1 + a_{22}v_2 \end{bmatrix}$$

# Computation of $\partial F / \partial x$ : Explicit Formulas

$$v_L = b_L + A_L v_{L-1} \quad \text{or simply} \quad w = b + Av.$$

The output  $w_i$  is not affected by  $b_j$  or  $A_{jk}$  if  $j \neq i$

Fully connected layer

Independent weights  $A_{jk}$

$$\frac{\partial w_i}{\partial b_j} = \delta_{ij} \quad \text{and} \quad \frac{\partial w_i}{\partial A_{jk}} = \delta_{ij} v_k$$

Example 
$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} + \begin{bmatrix} a_{11}v_1 + a_{12}v_2 \\ a_{21}v_1 + a_{22}v_2 \end{bmatrix}$$

$$\frac{\partial w_1}{\partial b_1} = \mathbf{1}, \quad \frac{\partial w_1}{\partial b_2} = \mathbf{0}, \quad \frac{\partial w_1}{\partial a_{11}} = \mathbf{v}_1, \quad \frac{\partial w_1}{\partial a_{12}} = \mathbf{v}_2, \quad \frac{\partial w_1}{\partial a_{21}} = \frac{\partial w_1}{\partial a_{22}} = \mathbf{0}.$$

# Backpropagation and the Chain Rule

---

$L(\mathbf{x})$  adds up all the losses  $\ell(\mathbf{w} - \mathbf{true}) = \ell(F(\mathbf{x}, \mathbf{v}) - \mathbf{true})$

# Backpropagation and the Chain Rule

---

$L(x)$  adds up all the losses  $\ell(w - \text{true}) = \ell(F(x, v) - \text{true})$

**The partial derivatives of  $L$  with respect to the weights  $x$  should be zero.**

# Backpropagation and the Chain Rule

$L(x)$  adds up all the losses  $\ell(w - \text{true}) = \ell(F(x, v) - \text{true})$

The partial derivatives of  $L$  with respect to the weights  $x$  should be zero.

**Chain rule**  $\frac{d}{dx}(F_3(F_2(F_1(x)))) = \left(\frac{dF_3}{dF_2}(F_2(F_1(x)))\right) \left(\frac{dF_2}{dF_1}(F_1(x))\right) \left(\frac{dF_1}{dx}(x)\right)$

# Backpropagation and the Chain Rule

$L(x)$  adds up all the losses  $\ell(w - \mathbf{true}) = \ell(F(x, v) - \mathbf{true})$

The partial derivatives of  $L$  with respect to the weights  $x$  should be zero.

**Chain rule**  $\frac{d}{dx}(F_3(F_2(F_1(x)))) = \left(\frac{dF_3}{dF_2}(F_2(F_1(x)))\right) \left(\frac{dF_2}{dF_1}(F_1(x))\right) \left(\frac{dF_1}{dx}(x)\right)$

What is the **multivariable** chain rule?

# Backpropagation and the Chain Rule

$L(x)$  adds up all the losses  $\ell(w - \mathbf{true}) = \ell(F(x, v) - \mathbf{true})$

The partial derivatives of  $L$  with respect to the weights  $x$  should be zero.

**Chain rule**  $\frac{d}{dx}(F_3(F_2(F_1(x)))) = \left(\frac{dF_3}{dF_2}(F_2(F_1(x)))\right) \left(\frac{dF_2}{dF_1}(F_1(x))\right) \left(\frac{dF_1}{dx}(x)\right)$

What is the **multivariable** chain rule?

Which order (**forward or backward** along the chain) is faster?

# Backward-mode $AD$ is faster for $M_1M_2\mathbf{w}$

$(M_1M_2)\mathbf{w}$  needs  $N^3+N^2$  multiplications     $M_1(M_2\mathbf{w})$  needs only  $N^2+N^2$

# Backward-mode $AD$ is faster for $M_1M_2\mathbf{w}$

$(M_1M_2)\mathbf{w}$  needs  $N^3+N^2$  multiplications     $M_1(M_2\mathbf{w})$  needs only  $N^2+N^2$

**Forward**  $((M_1M_2)M_3)\dots M_L\mathbf{w}$  needs  $(L-1)N^3 + N^2$

# Backward-mode $AD$ is faster for $M_1M_2\mathbf{w}$

$(M_1M_2)\mathbf{w}$  needs  $N^3+N^2$  multiplications     $M_1(M_2\mathbf{w})$  needs only  $N^2+N^2$

**Forward**  $((M_1M_2)M_3)\dots M_L\mathbf{w}$  needs  $(L-1)N^3+N^2$

**Backward**  $M_1(M_2(\dots(M_L\mathbf{w})))$  needs  $LN^2$

# The Multivariable Chain Rule

$$\frac{\partial \mathbf{w}}{\partial \mathbf{v}} = \begin{bmatrix} \frac{\partial w_1}{\partial v_1} & \cdots & \frac{\partial w_1}{\partial v_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial w_p}{\partial v_1} & \cdots & \frac{\partial w_p}{\partial v_n} \end{bmatrix}$$

$$\frac{\partial \mathbf{v}}{\partial \mathbf{u}} = \begin{bmatrix} \frac{\partial v_1}{\partial u_1} & \cdots & \frac{\partial v_1}{\partial u_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial v_n}{\partial u_1} & \cdots & \frac{\partial v_n}{\partial u_m} \end{bmatrix}$$

# The Multivariable Chain Rule

$$\frac{\partial w}{\partial v} = \begin{bmatrix} \frac{\partial w_1}{\partial v_1} & \cdots & \frac{\partial w_1}{\partial v_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial w_p}{\partial v_1} & \cdots & \frac{\partial w_p}{\partial v_n} \end{bmatrix} \quad \frac{\partial v}{\partial u} = \begin{bmatrix} \frac{\partial v_1}{\partial u_1} & \cdots & \frac{\partial v_1}{\partial u_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial v_n}{\partial u_1} & \cdots & \frac{\partial v_n}{\partial u_m} \end{bmatrix}$$

$$\frac{\partial w_i}{\partial u_k} = \frac{\partial w_i}{\partial v_1} \frac{\partial v_1}{\partial u_k} + \cdots + \frac{\partial w_i}{\partial v_n} \frac{\partial v_n}{\partial u_k} = \left( \frac{\partial w_i}{\partial v_1}, \dots, \frac{\partial w_i}{\partial v_n} \right) \cdot \left( \frac{\partial v_1}{\partial u_k}, \dots, \frac{\partial v_n}{\partial u_k} \right)$$

**Multivariable chain rule: Multiply matrices!**  $\frac{\partial w}{\partial u} = \left( \frac{\partial w}{\partial v} \right) \left( \frac{\partial v}{\partial u} \right)$

# Hyperparameters: The Fateful Decisions

---

The words **learning rate** are often used in place of stepsize

# Hyperparameters: The Fateful Decisions

---

The words **learning rate** are often used in place of stepsize

$s_k$  **is too small** Then gradient descent takes too long to minimize  $L(\mathbf{x})$

# Hyperparameters: The Fateful Decisions

The words **learning rate** are often used in place of stepsize

$s_k$  **is too small** Then gradient descent takes too long to minimize  $L(\boldsymbol{x})$

$s_k$  **is too large** Overshooting the best choice  $\boldsymbol{x}_{k+1}$  in the descent direction

# Hyperparameters: The Fateful Decisions

The words **learning rate** are often used in place of stepsize

$s_k$  **is too small** Then gradient descent takes too long to minimize  $L(\mathbf{x})$

$s_k$  **is too large** Overshooting the best choice  $\mathbf{x}_{k+1}$  in the descent direction

**Cross-validation** Divide the available data into  $K$  subsets

# Regularization = Weight decay : $\ell^2$ or $\ell^1$

---

Small  $\lambda$  : increase the variance of the error (overfitting)

# Regularization = Weight decay: $\ell^2$ or $\ell^1$

Small  $\lambda$ : increase the variance of the error (overfitting)

Large  $\lambda$ : increase the bias (underfitting),  $\|\mathbf{b} - A\mathbf{x}\|^2$  is less important

# Regularization = Weight decay: $\ell^2$ or $\ell^1$

Small  $\lambda$ : increase the variance of the error (overfitting)

Large  $\lambda$ : increase the bias (underfitting),  $\|\mathbf{b} - A\mathbf{x}\|^2$  is less important

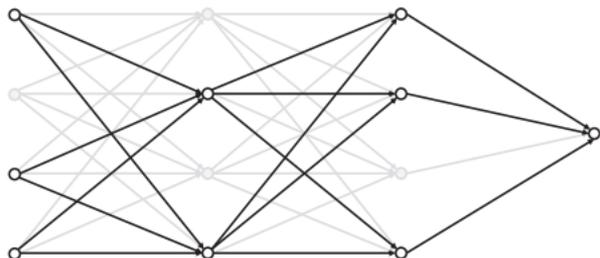
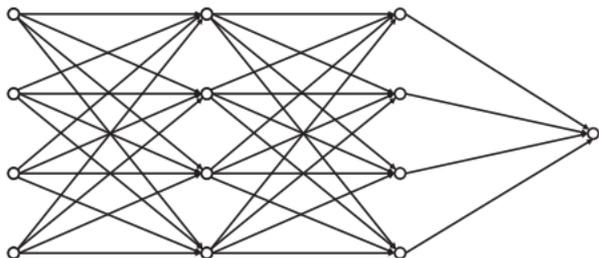
Deep learning with many extra weights and good hyperparameters will find solutions that generalize, without penalty

# Regularization = Weight decay : $\ell^2$ or $\ell^1$

Small  $\lambda$  : increase the variance of the error (overfitting)

Large  $\lambda$  : increase the bias (underfitting),  $\|\mathbf{b} - A\mathbf{x}\|^2$  is less important

Deep learning with many extra weights and good hyperparameters will find solutions that generalize, without penalty



# Softmax Outputs for Multiclass Networks

$$\text{Softmax} \quad p_j = \frac{1}{S} e^{w_j} \quad \text{where} \quad S = \sum_{k=1}^n e^{w_k}$$

# Softmax Outputs for Multiclass Networks

$$\text{Softmax} \quad p_j = \frac{1}{S} e^{w_j} \quad \text{where} \quad S = \sum_{k=1}^n e^{w_k}$$

Softmax produces the probabilities in **teachyourmachine.com**

# Softmax Outputs for Multiclass Networks

$$\text{Softmax} \quad p_j = \frac{1}{S} e^{w_j} \quad \text{where} \quad S = \sum_{k=1}^n e^{w_k}$$

Softmax produces the probabilities in [teachyourmachine.com](http://teachyourmachine.com)

**The World Championship at the Game of Go**

# Softmax Outputs for Multiclass Networks

$$\text{Softmax} \quad p_j = \frac{1}{S} e^{w_j} \quad \text{where} \quad S = \sum_{k=1}^n e^{w_k}$$

Softmax produces the probabilities in **teachyourmachine.com**

**The World Championship at the Game of Go**

**Residual Networks (ResNets)** “*skip connections*”

# Neural Nets Give Universal Approximation

**If  $f(v)$  is continuous there exists  $x$  so that  $|F(x, v) - f(v)| < \epsilon$  for all  $v$**

# Neural Nets Give Universal Approximation

If  $f(v)$  is continuous there exists  $x$  so that  $|F(x, v) - f(v)| < \epsilon$  for all  $v$

**Accuracy of approximation to  $f$**        $\min_x \|F(x, v) - f(v)\| \leq C\|f\|_S$

# Neural Nets Give Universal Approximation

If  $f(v)$  is continuous there exists  $x$  so that  $|F(x, v) - f(v)| < \epsilon$  for all  $v$

Accuracy of approximation to  $f$   $\min_x \|F(x, v) - f(v)\| \leq C\|f\|_S$

Deep networks give **closer approximation** than splines or shallow nets

# Counting Flat Pieces in the Graph

**Theorem** For  $\mathbf{v}$  in  $\mathbf{R}^m$ , suppose the graph of  $F(\mathbf{v})$  has folds along  $N$  hyperplanes  $H_1, \dots, H_N$ . Those come from  $N$  linear equations  $\mathbf{a}_i^T \mathbf{v} + b_i = 0$ , in other words ReLU at  $N$  neurons.  $F$  has  $r(N, m)$  linear pieces:

$$r(N, m) = \sum_{i=0}^m \binom{N}{i} = \binom{N}{\mathbf{0}} + \binom{N}{\mathbf{1}} + \dots + \binom{N}{\mathbf{m}}$$

# Counting Flat Pieces in the Graph

**Theorem** For  $\mathbf{v}$  in  $\mathbf{R}^m$ , suppose the graph of  $F(\mathbf{v})$  has folds along  $N$  hyperplanes  $H_1, \dots, H_N$ . Those come from  $N$  linear equations  $\mathbf{a}_i^T \mathbf{v} + b_i = 0$ , in other words ReLU at  $N$  neurons.  $F$  has  $r(N, m)$  linear pieces:

$$r(N, m) = \sum_{i=0}^m \binom{N}{i} = \binom{N}{0} + \binom{N}{1} + \dots + \binom{N}{m}$$

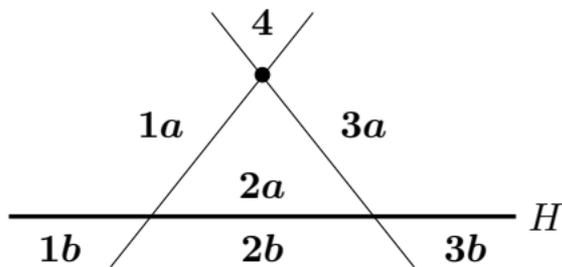
$$r(N, m) = r(N - 1, m) + r(N - 1, m - 1)$$

# Counting Flat Pieces in the Graph

**Theorem** For  $v$  in  $\mathbf{R}^m$ , suppose the graph of  $F(v)$  has folds along  $N$  hyperplanes  $H_1, \dots, H_N$ . Those come from  $N$  linear equations  $\mathbf{a}_i^T v + b_i = 0$ , in other words ReLU at  $N$  neurons.  $F$  has  $r(N, m)$  linear pieces:

$$r(N, m) = \sum_{i=0}^m \binom{N}{i} = \binom{N}{0} + \binom{N}{1} + \dots + \binom{N}{m}$$

$$r(N, m) = r(N - 1, m) + r(N - 1, m - 1)$$



Start with 2 planes

$$\leftarrow r(2, 2) = 4$$

Add new plane  $H$

$$\leftarrow r(2, 1) = 3$$

# Continuous Piecewise Linear Function

---

How many linear pieces with more layers?

# Continuous Piecewise Linear Function

---

How many linear pieces with more layers ?

Now ReLU is folding **piecewise linear** functions

# Continuous Piecewise Linear Function

How many linear pieces with more layers ?

Now ReLU is folding **piecewise linear** functions

Hanin-Rolnick: Still  $r(N, m) \approx cN^m$  pieces from  $N$  neurons

MIT OpenCourseWare

<https://ocw.mit.edu>

## 18.085 Computational Science and Engineering I

Summer 2020

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.