

# **RBE474X/595-B01-ST: Deep Learning For Perception**

**Class 2: Multi Layer Perceptrons And Backpropagation**

**Prof. Wei Xiao**



**Shreyas Devdatta Khobragade**

*he/him/his*

Office Hour Location: UH243 Curtain Space

Office Hours: Tuesday 2:00pm-3:00pm



[https://piazza.com  
/wpi/fall2025/rbe4  
74x595](https://piazza.com/wpi/fall2025/rbe474x595)

# Please Ask Questions!

CAN UNICORNS ROAST  
MARSHMALLOWS?

NO.



**Today's lecture is Math heavy! Slow me down if needed!**

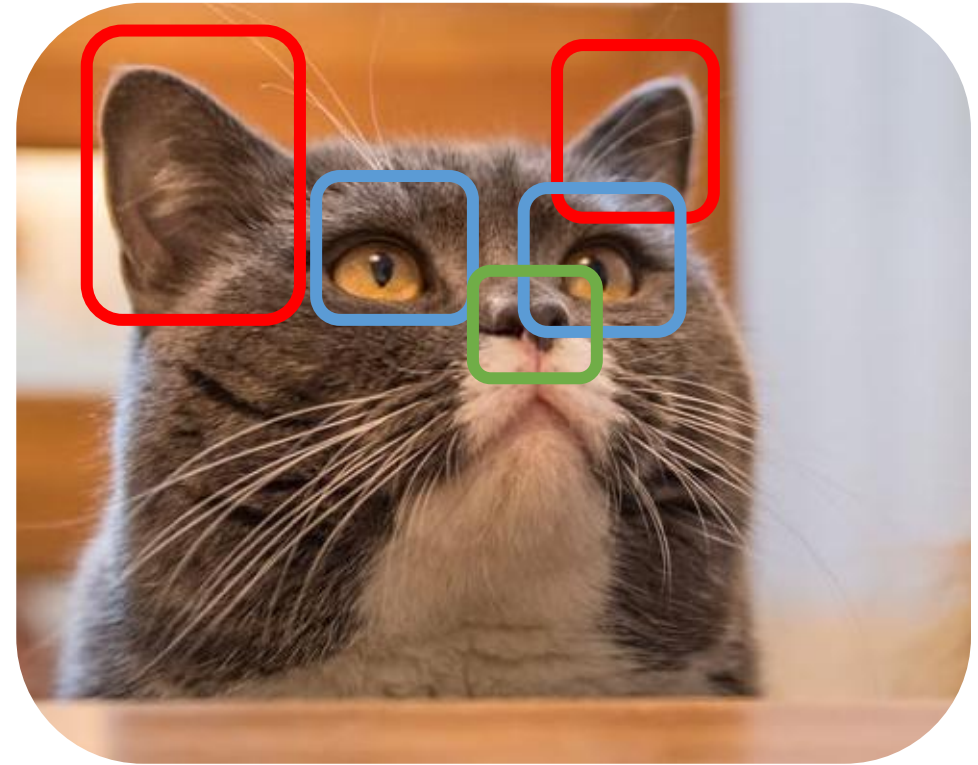
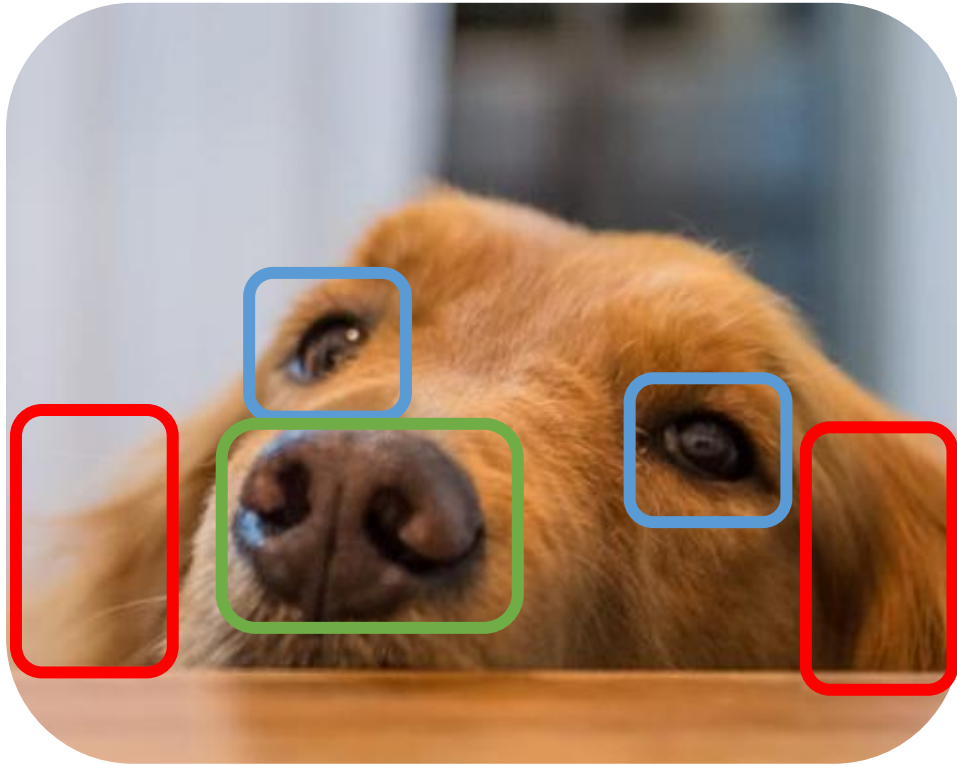


# Dog or Cat

## Image Classification

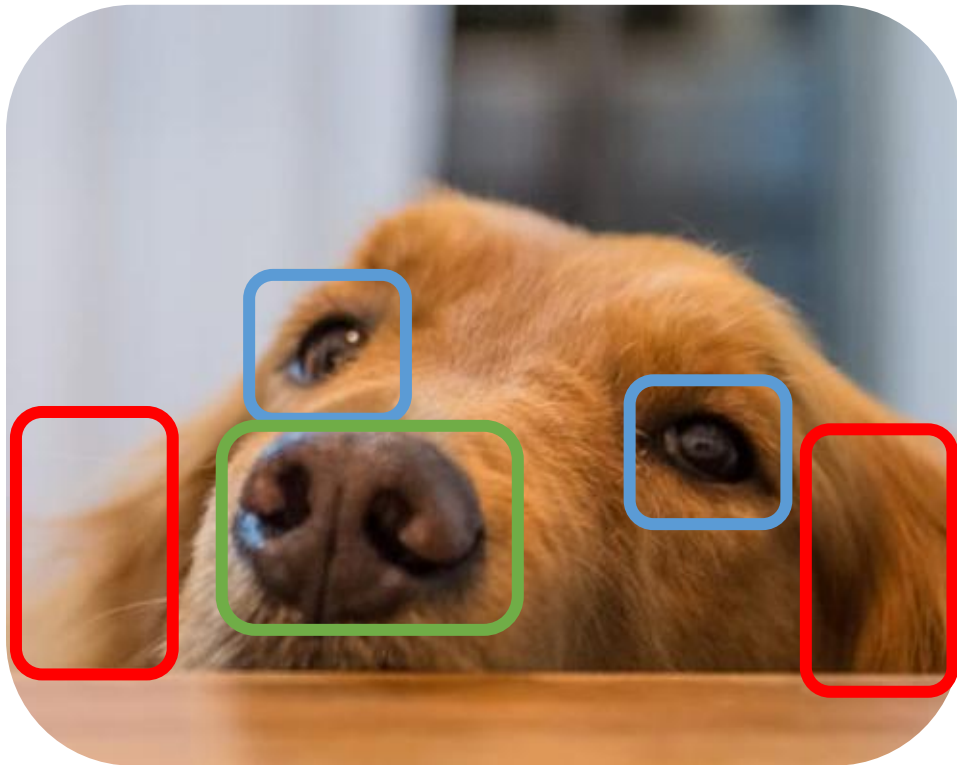
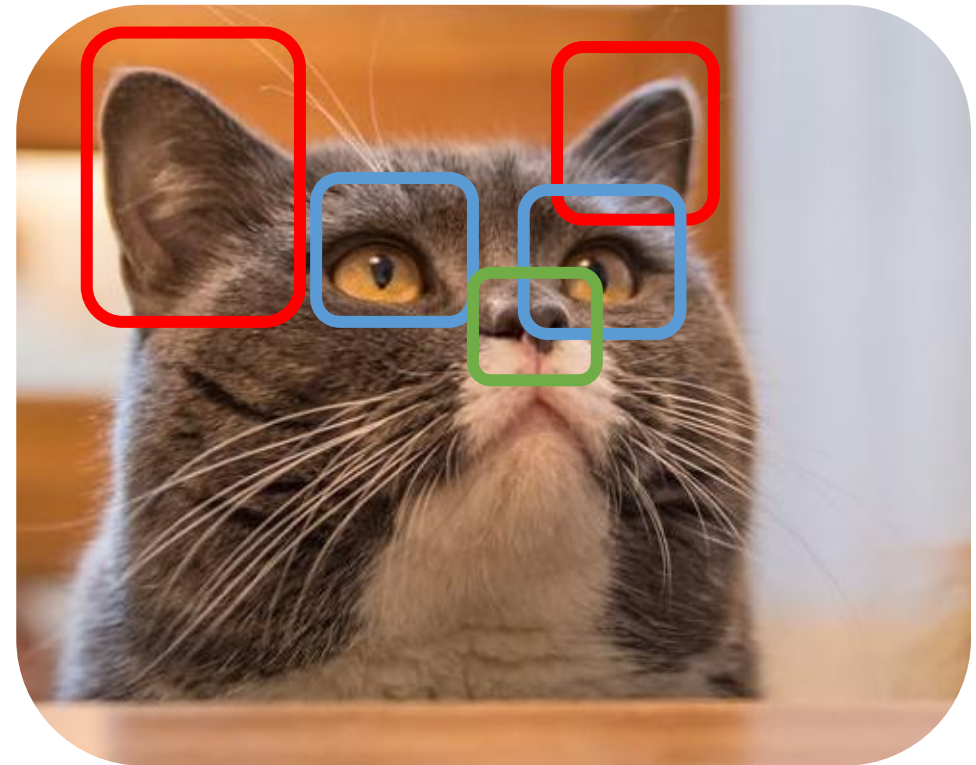


# How Do You Tell If It's A Dog or a Cat?

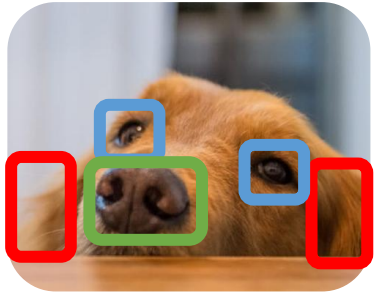


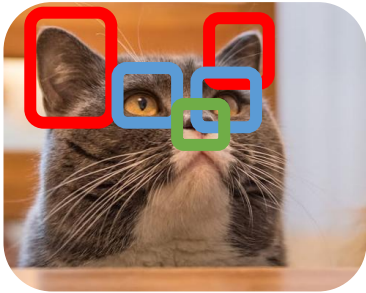
# Features

$$X = \begin{bmatrix} \text{Non - pointy ears} \\ \text{Round and large nose} \\ \text{Eyes without pointy pupil} \end{bmatrix}$$


$$\begin{bmatrix} 0.21 \\ 0.99 \\ 0.86 \end{bmatrix}$$

$$\begin{bmatrix} 0.79 \\ 0.01 \\ 0.14 \end{bmatrix}$$

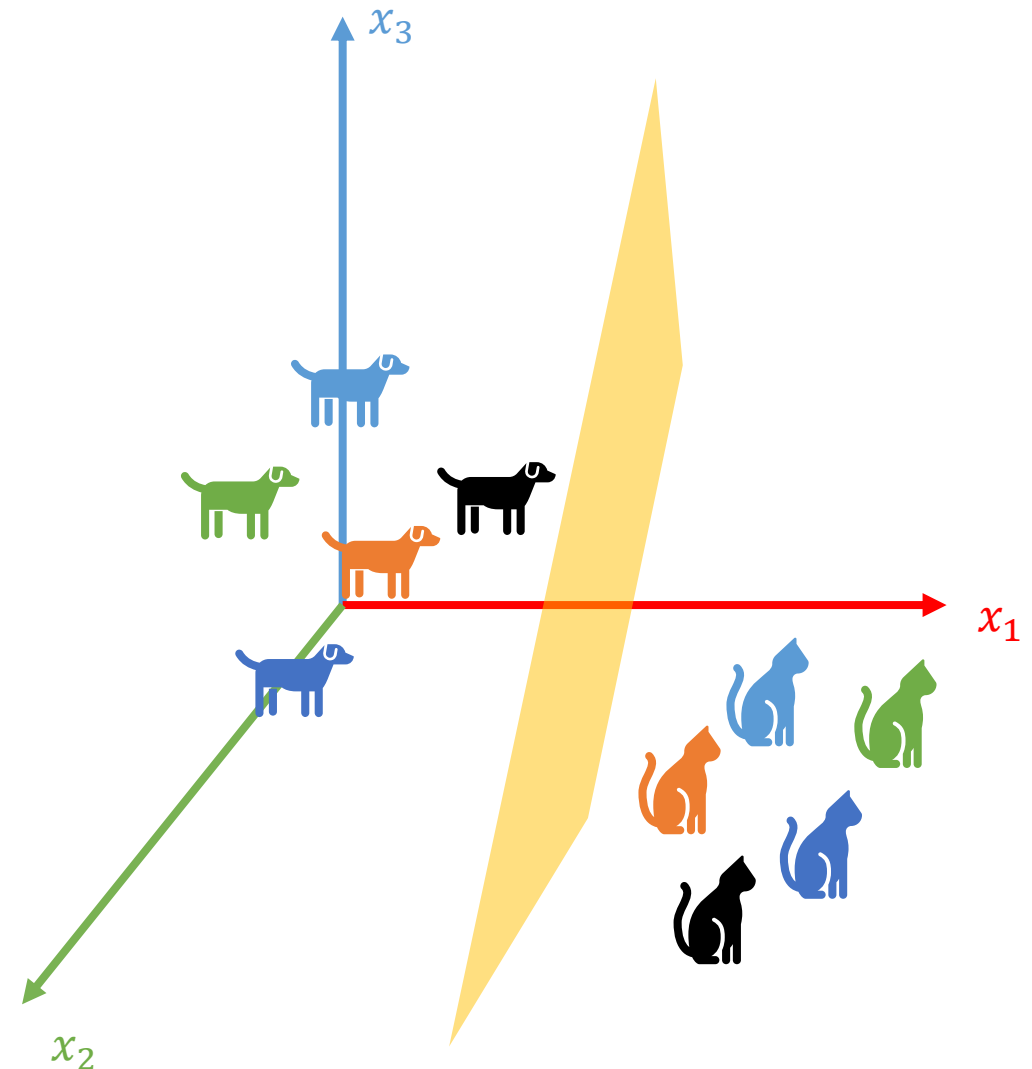
# Features



$$\begin{bmatrix} 0.21 \\ 0.99 \\ 0.86 \end{bmatrix}$$


$$\begin{bmatrix} 0.79 \\ 0.01 \\ 0.14 \end{bmatrix}$$

$$X = \begin{bmatrix} \text{Non - pointy ears} \\ \text{Round and large nose} \\ \text{Eyes without pointy pupil} \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$





# Linear Regression

$$Y = W_0^T X + B$$

$$\text{or } Y = W^T \begin{bmatrix} X \\ 1 \end{bmatrix}, W = [W_0^T \ B]$$

Given  $\{X, Y\}_i$  find  $W$

Mathematically,

$$\operatorname{argmin}_W \mathbb{E} \left( \|\hat{Y} - \tilde{Y}\|_2 \right)$$

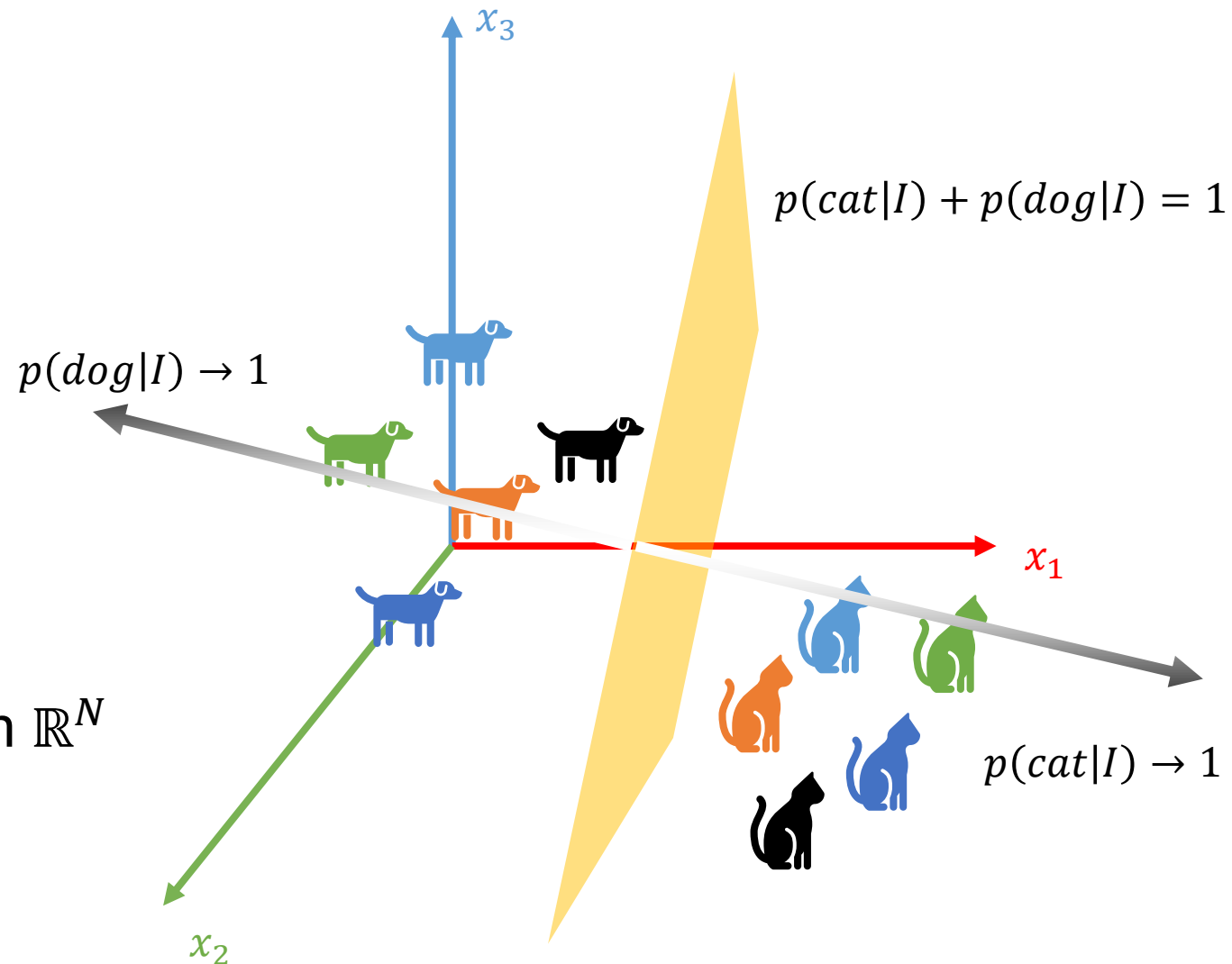
$$\text{or } \operatorname{argmin}_W \mathbb{E} \left( \|\hat{Y} - f(X, W)\|_2 \right)$$

We are fitting a linear hyperplane in  $\mathbb{R}^N$

How do you get  $p$  from  $Y$ ?

You need some function!

$$p = g(\tilde{Y})$$





# Linear Regression

$$Y = W^T \begin{bmatrix} X \\ 1 \end{bmatrix}$$

$$Y = AW$$

$$W = A^{-1}Y$$

When is this valid?

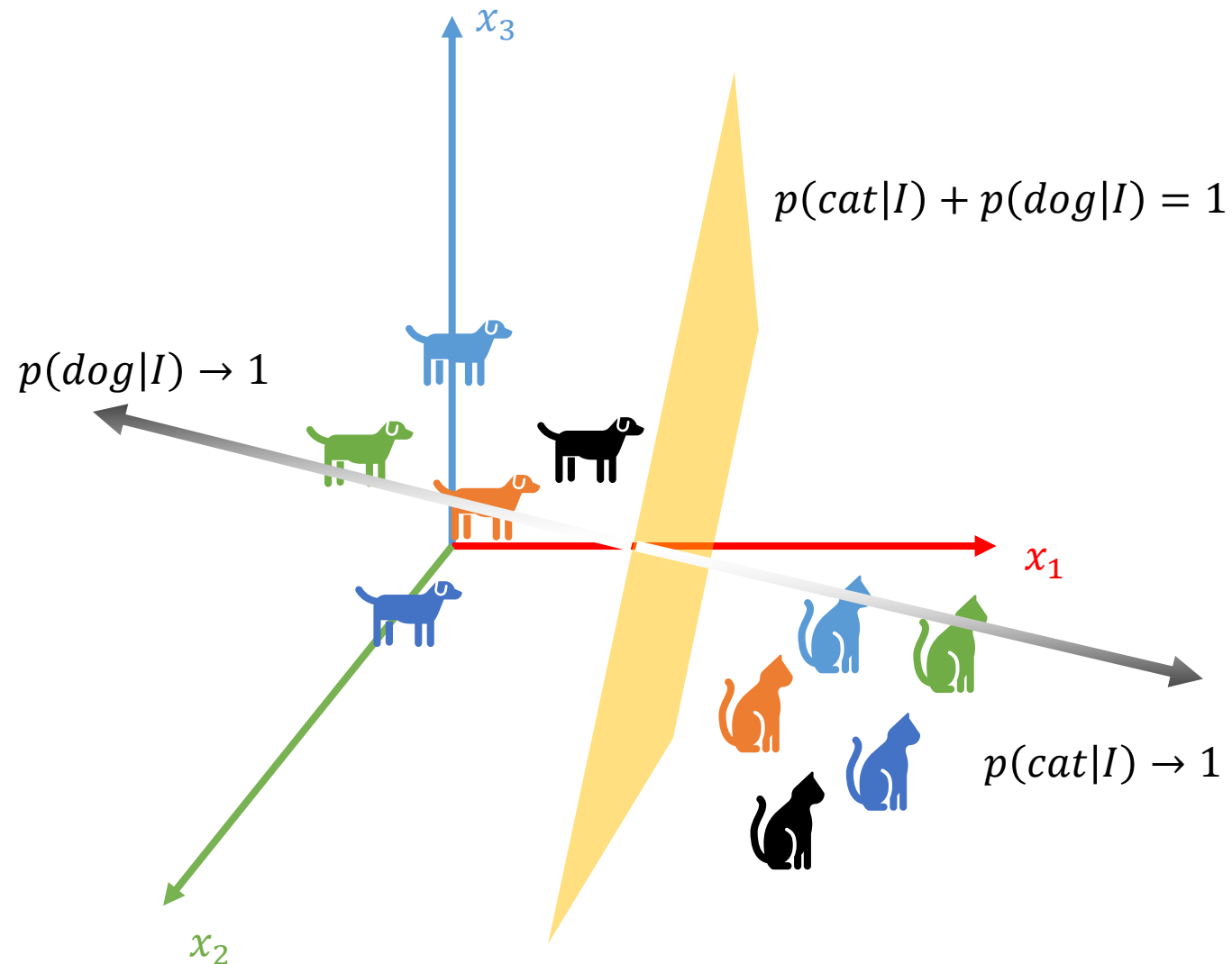
$A$  is square and invertible!

What if  $A$  is not square?

$$A^{-1} \approx A^\dagger$$

How do you get the pseudo-inverse?

SVD!



# Singular Value Decomposition

Dimensionality of Data

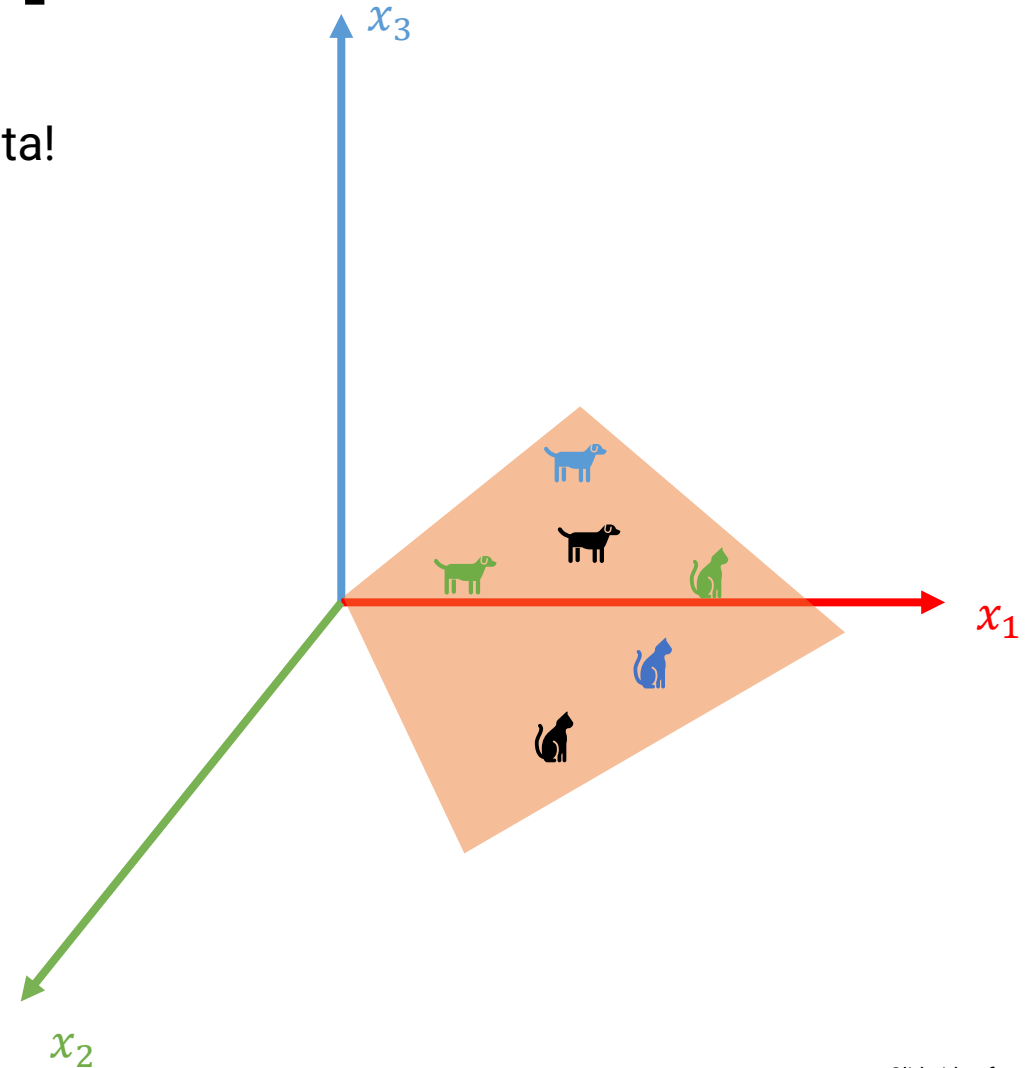
Min. num. independent dimensions to represent **all** data!

How do we know this mathematically?

Rank of  $A$ !

In our example, what is the dimensionality?

2! But it is represented in 3D.



Slide idea from Stanford's CS426

# Singular Value Decomposition

$\Sigma$  can also be shown as  $D$

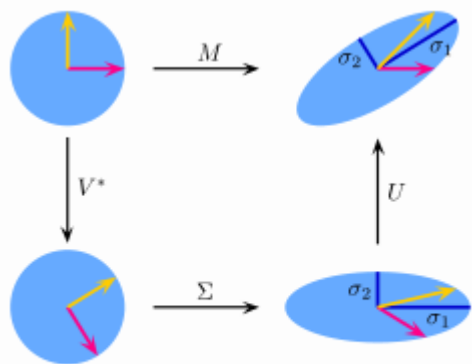
$$A = U \Sigma V^T$$

$\begin{matrix} \nearrow & \downarrow & \downarrow \\ M \times N & M \times M & N \times N \\ & \downarrow & \\ & M \times N & \end{matrix}$

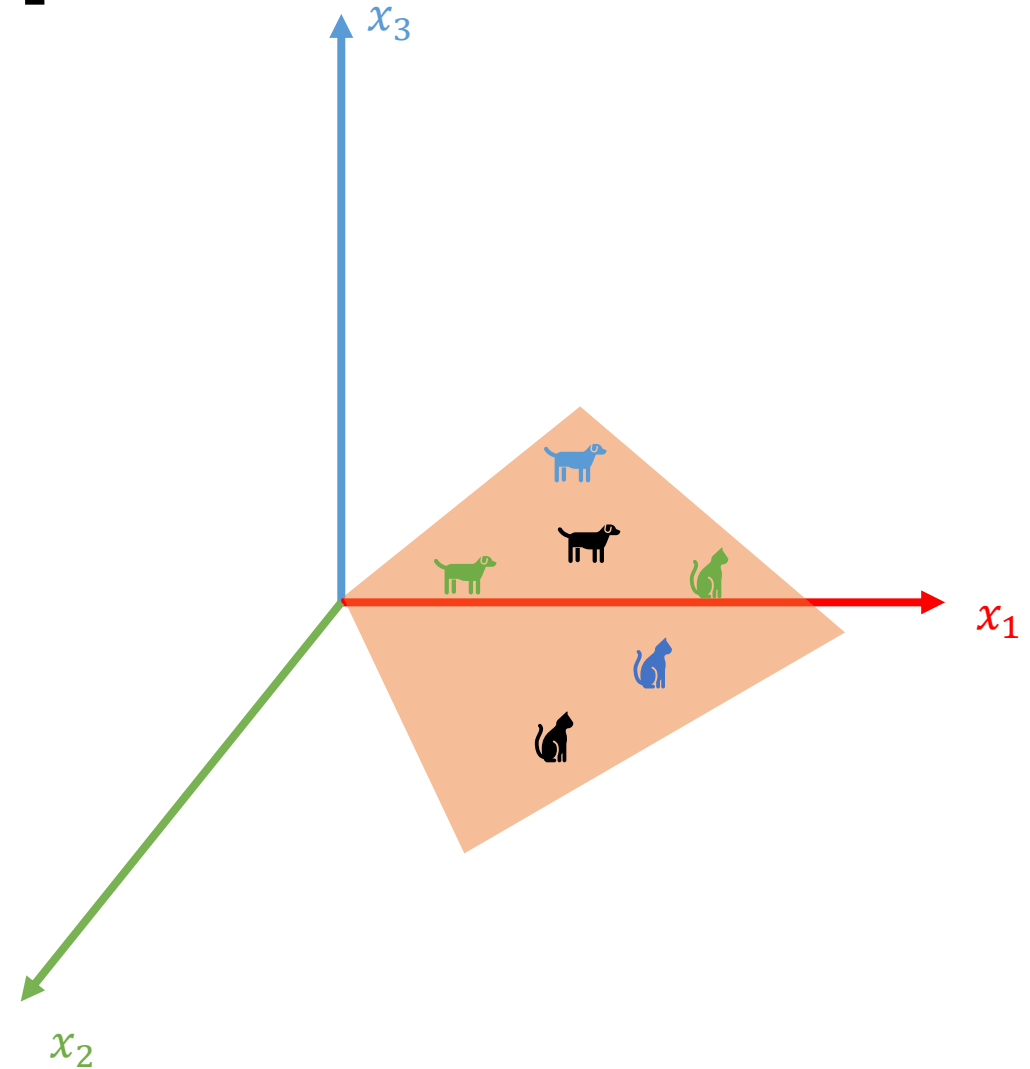
Shows the spread of data!

$$UU^T = \mathbb{I}_M$$

$$VV^T = \mathbb{I}_N$$



$$M = U \cdot \Sigma \cdot V^*$$



# SVD

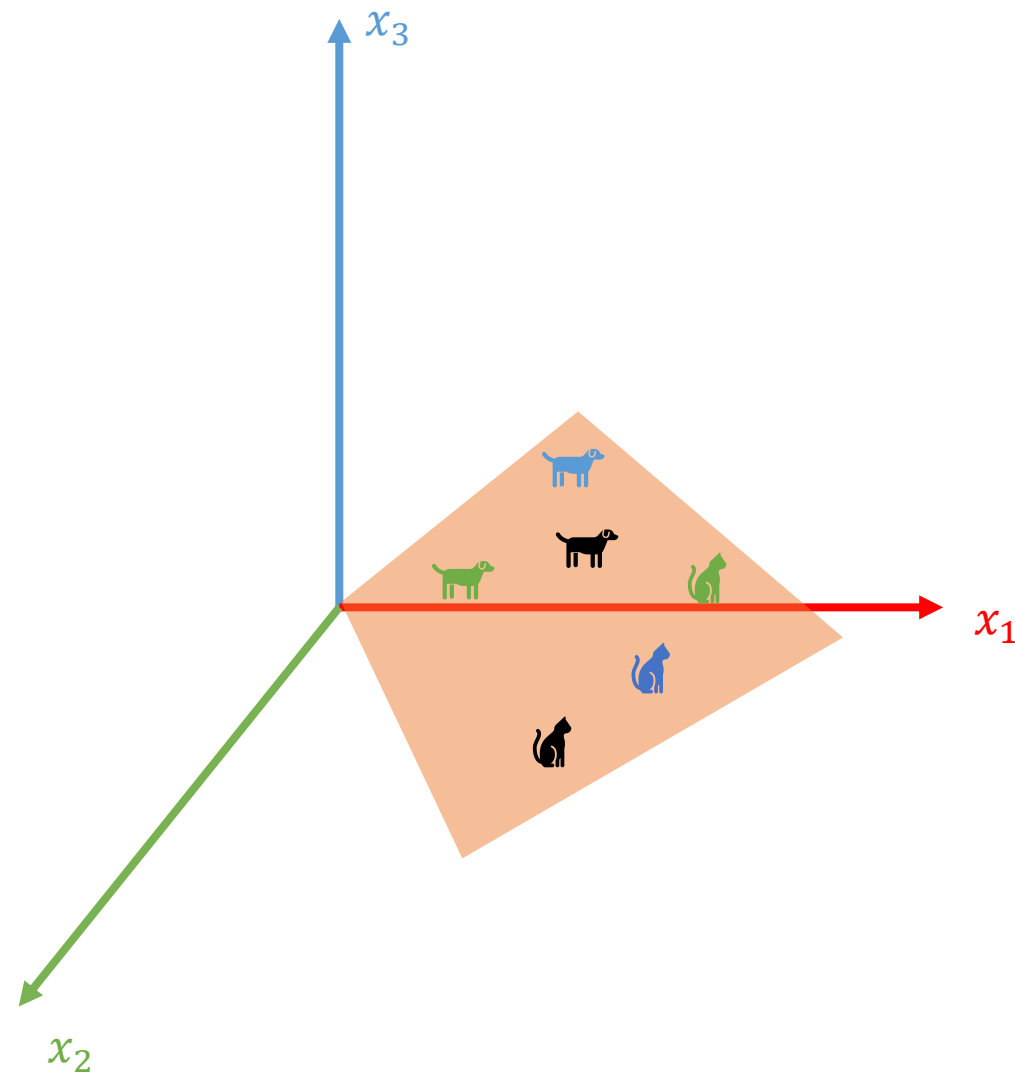
$$A \mathbf{x} \approx \mathbf{b}$$

$$A_{M \times N} = U_{M \times M} D_{M \times N} V^T_{N \times N}$$

$U$  Orthogonal matrix  $UU^T = \mathbb{I}_M$   
 $D$  Diagonal matrix  
 $V^T$  Orthogonal matrix

What do you think num. red squares is?

The **true** dimensionality!  
**The number of actual min. dim. needed to represent data!**  
 Also called Num. Principal Components in PCA





# Nullspace

$$\begin{matrix} \boxed{A} & \boxed{\mathbf{x}} & = & \boxed{\mathbf{0}} \\ M \times N & N \times 1 & & M \times 1 \end{matrix}$$

$\mathbf{x} = \text{null}(A)$ , the set of all vectors  $\mathbf{x}$  that satisfy the above equation (the rank of  $A$  is  $\rho \leq \min(M, N)$ )

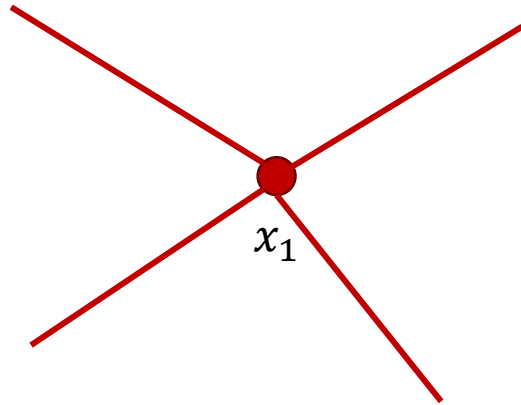
$$\begin{matrix} \boxed{A} & = & \boxed{U} & \boxed{D} & \boxed{V^T} \\ M \times N & & M \times M & M \times N & N \times N \\ & & \text{Orthogonal matrix} & \text{Diagonal matrix} & \text{Orthogonal matrix} \end{matrix}$$

Diagram details: The diagonal matrix  $D$  is shown with a red box containing  $\rho$  diagonal elements and a white box containing  $N - \rho$  zeros. The orthogonal matrix  $V^T$  is shown with a yellow box containing  $\rho$  rows and a white box containing  $N - \rho$  rows.

$$V_{\rho+1:N} = \text{null}(A)$$

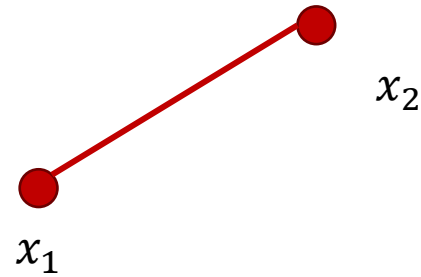
If  $\rho = N$ , the matrix has full rank, and the null space is unique (i.e., 0)

# Line Fitting 101



Under constrained problem:  $\infty$  Solutions

# Line Fitting 101



$$ax + by + c = 0$$

$$ax_1 + by_1 + c = 0$$

$$ax_2 + by_2 + c = 0$$

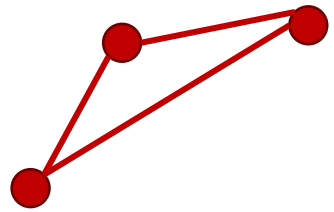
Form  $A\mathbf{x} = \mathbf{0}, \mathbf{x} = (\mathbf{a}, \mathbf{b}, \mathbf{c})$

Perfectly constrained problem: Unique Solution

Can be written as  $\operatorname{argmin}_{\mathbf{x}=(a,b,c)} \|A\mathbf{x}\|_2^2 \text{ s.t. } \|\mathbf{x}\|_2 = 1$

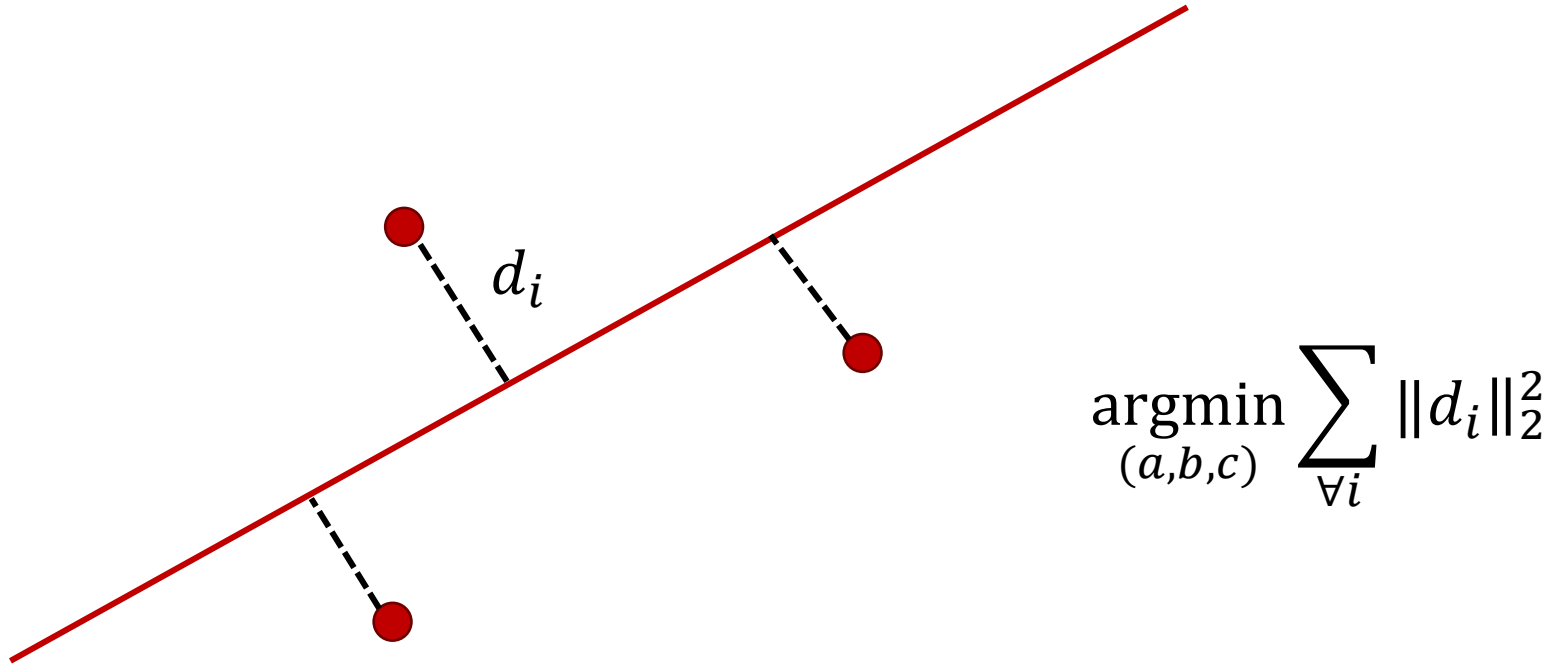
↓  
Avoids trivial solution

# Line Fitting 101





# Line Fitting 101



$$A\mathbf{x} = \mathbf{b}$$

Instead of  $ax + by + c = 0$   
Let's use  $y = mx + d$

$$\begin{array}{|c|} \hline A \\ \hline \end{array} \begin{array}{|c|} \hline \mathbf{x} \\ \hline \end{array} \approx \begin{array}{|c|} \hline \mathbf{b} \\ \hline \end{array}$$

How do we solve this?

Use a loss function (metric) that needs to be minimized (optimized)

Let's minimize  $\arg\min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|_2^2$

How do you obtain closed form solution?

Differentiate and set to zero!

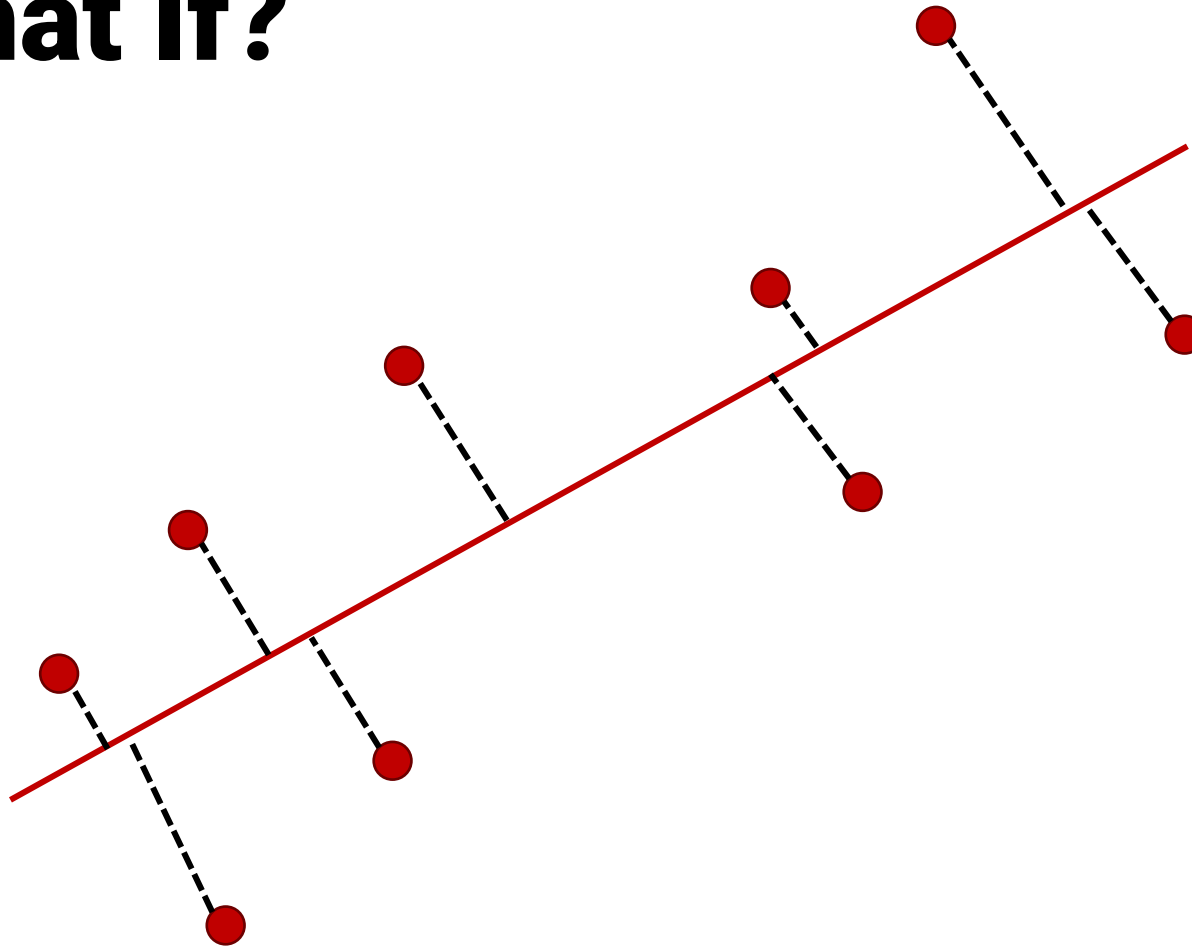
Solved using  $\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$

Called **Ordinary Least Squares (OLS)**

$$\mathbf{x} = A \backslash \mathbf{b};$$



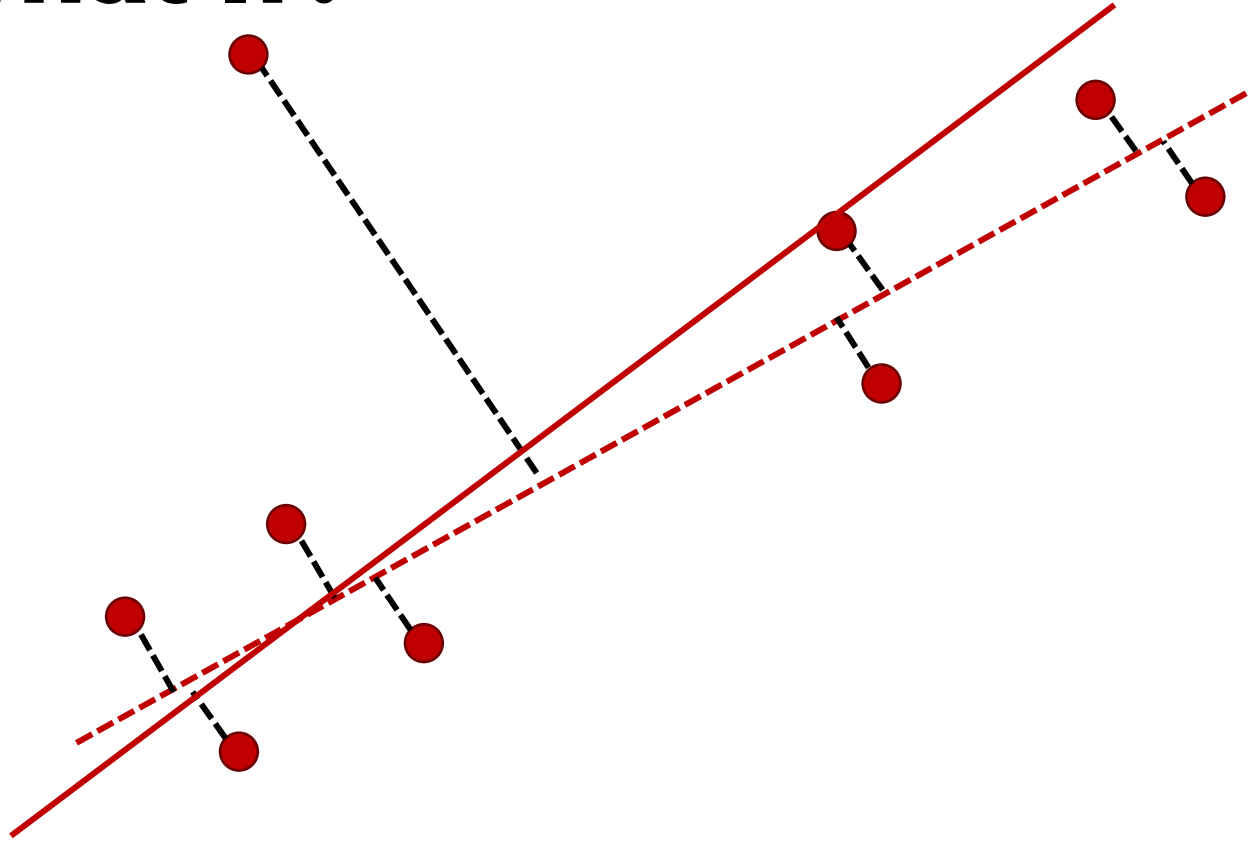
# What If?



$$\operatorname{argmin}_{(a,b,c)} \sum_{\forall i} \|d_i\|_2^2$$

Same as before!

# What If?



$$\operatorname{argmin}_{(a,b,c)} \sum_{\forall i} \|d_i\|_2^2$$

Outliers!

You can regularize!

$$\operatorname{argmin}_{(a,b,c)} \sum_{\forall i} \|d_i\|_2^2 + \lambda \left\| \begin{bmatrix} a \\ b \\ c \end{bmatrix} \right\|_2^2$$

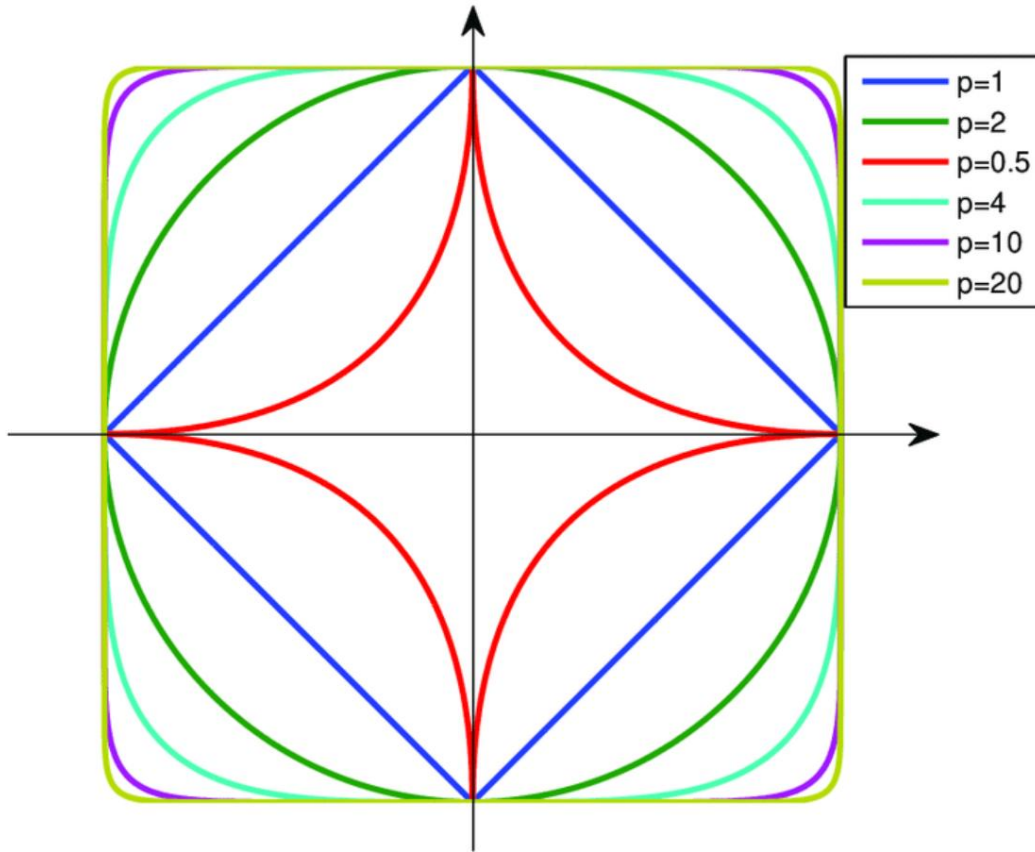
Penalize for large values in model!  
(This penalty discourages the model from assigning large coefficients to features)

Penalty or Regularization can be  $l_p$  norm

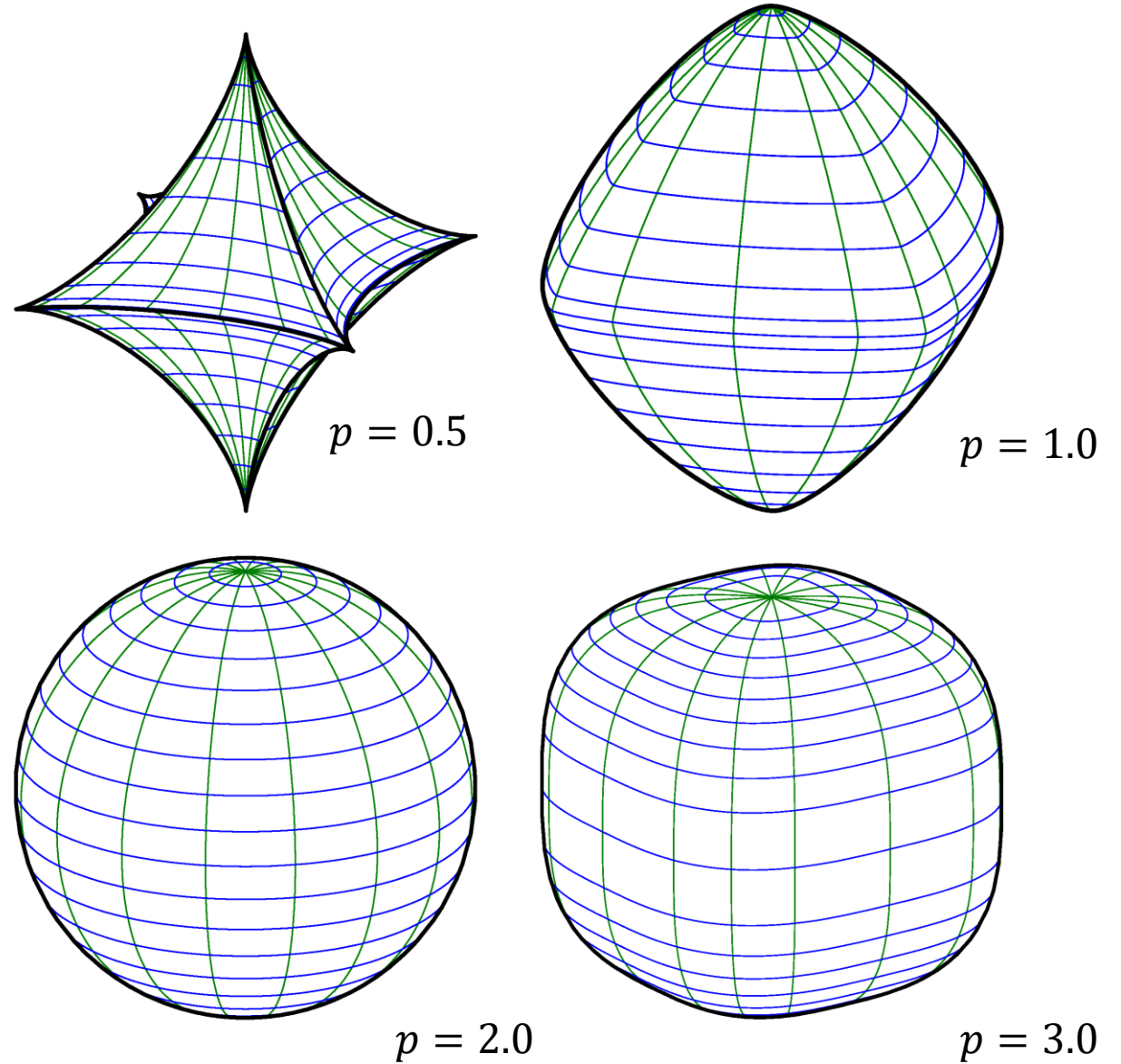
$$l_p \text{ norm: } \|\mathbf{x}\|_p = \left( \sum_{\forall i} x_i^p \right)^{\frac{1}{p}}$$



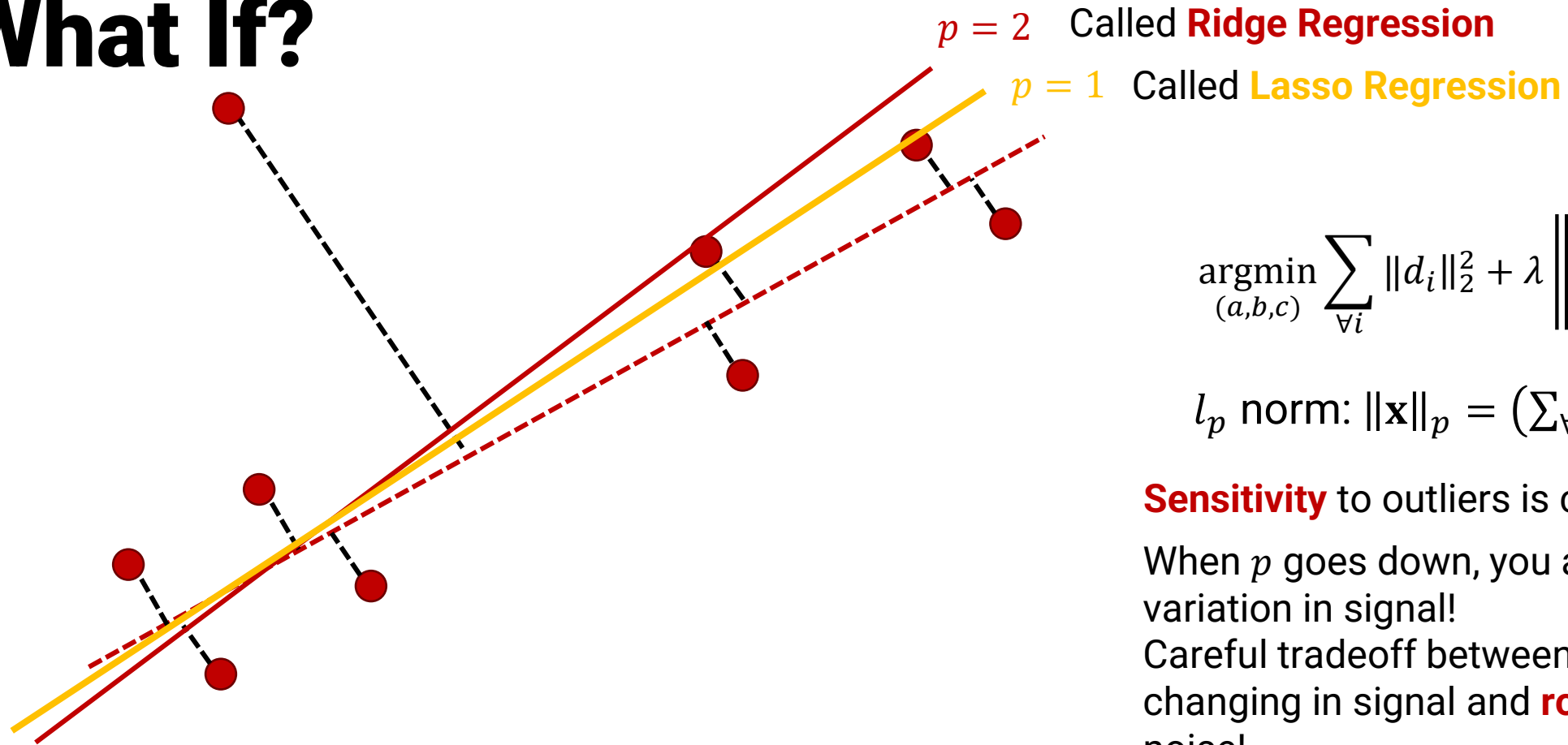
# World Of Norms!



Help alleviate different kinds of noise!



# What If?



$$\operatorname{argmin}_{(a,b,c)} \sum_{\forall i} \|d_i\|_2^2 + \lambda \left\| \begin{bmatrix} a \\ b \\ c \end{bmatrix} \right\|_p$$

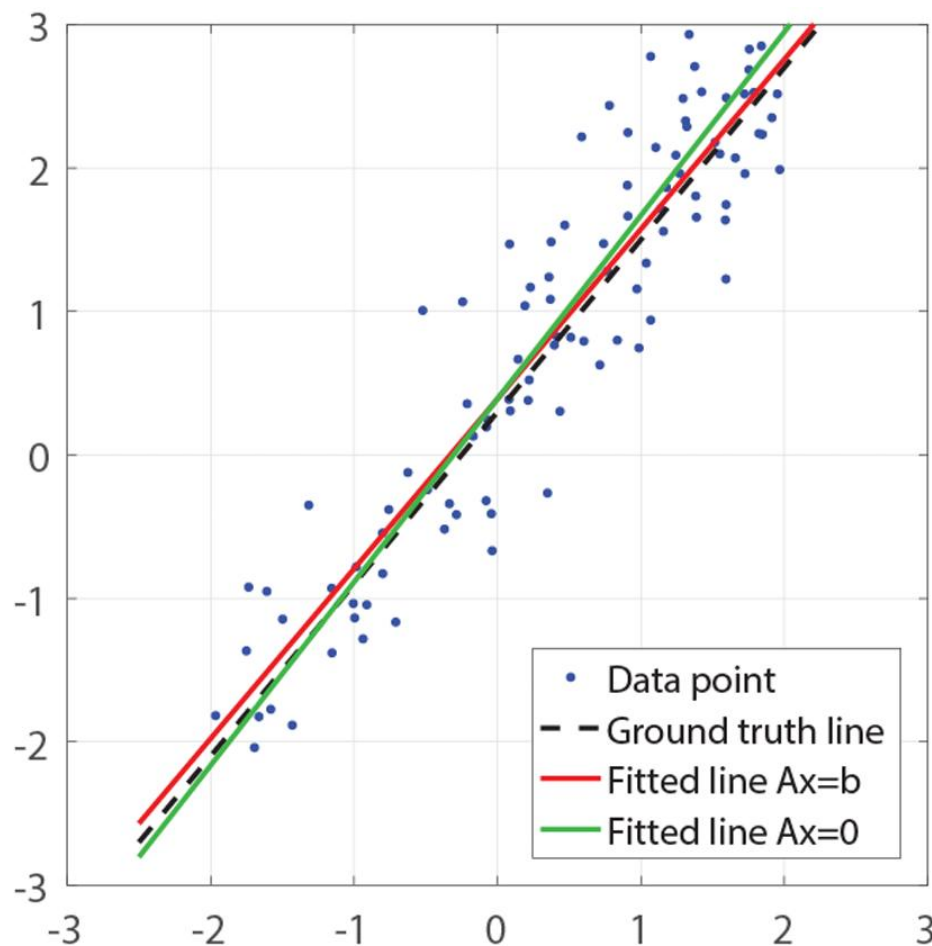
$$l_p \text{ norm: } \|\mathbf{x}\|_p = \left( \sum_{\forall i} x_i^p \right)^{\frac{1}{p}}$$

**Sensitivity** to outliers is controlled by  $p$ !

When  $p$  goes down, you also reject variation in signal!

Careful tradeoff between **sensitivity** to changing in signal and **robustness** to noise!

# $A\mathbf{x} = \mathbf{0}$ Vs $A\mathbf{x} = \mathbf{b}$



$$\begin{aligned} au_1 + bv_1 + c &\approx 0 \\ au_2 + bv_2 + c &\approx 0 \\ &\vdots \\ au_n + bv_n + c &\approx 0 \end{aligned}$$

$$\begin{aligned} v_1 &\approx mu_1 + d \\ v_2 &\approx mu_2 + d \\ &\vdots \\ v_n &\approx mu_n + d \end{aligned}$$

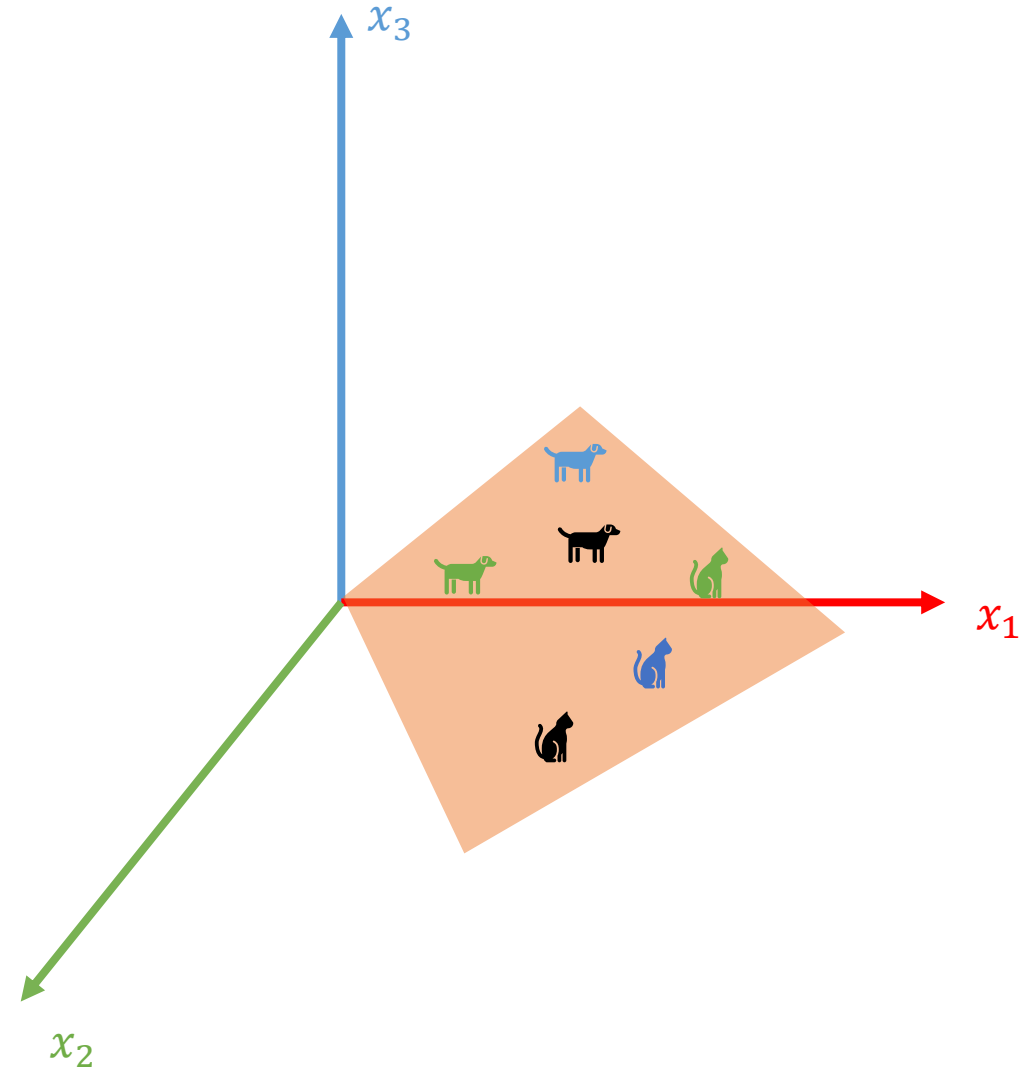
$$\begin{bmatrix} u_1 & v_1 & 1 \\ u_2 & v_2 & 1 \\ \vdots & \vdots & \vdots \\ u_n & v_n & 1 \end{bmatrix} \mathbf{A} \begin{bmatrix} a \\ b \\ c \end{bmatrix} \mathbf{x} \approx \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} u_1 & 1 \\ u_2 & 1 \\ \vdots & \vdots \\ u_n & 1 \end{bmatrix} \mathbf{A} \begin{bmatrix} m \\ d \end{bmatrix} \mathbf{x} \approx \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \mathbf{b}$$

$$\begin{array}{c} \mathbf{Ax} = \mathbf{0} \\ \hline \begin{array}{ccc} \text{No sol.} & \text{Unique sol.} & \infty \text{ sols.} \\ \text{Inconsistent} & \text{Consistent} & \end{array} \end{array}$$

$$\begin{array}{c} \mathbf{Ax} = \mathbf{b} \\ \hline \begin{array}{cc} \text{Unique sol.} & \infty \text{ sols.} \\ \text{Trivial sol.} & \text{Non-trivial sol.} \\ \mathbf{x} = \mathbf{0} & \mathbf{x} \neq \mathbf{0} \end{array} \end{array}$$

# Noise In SVD

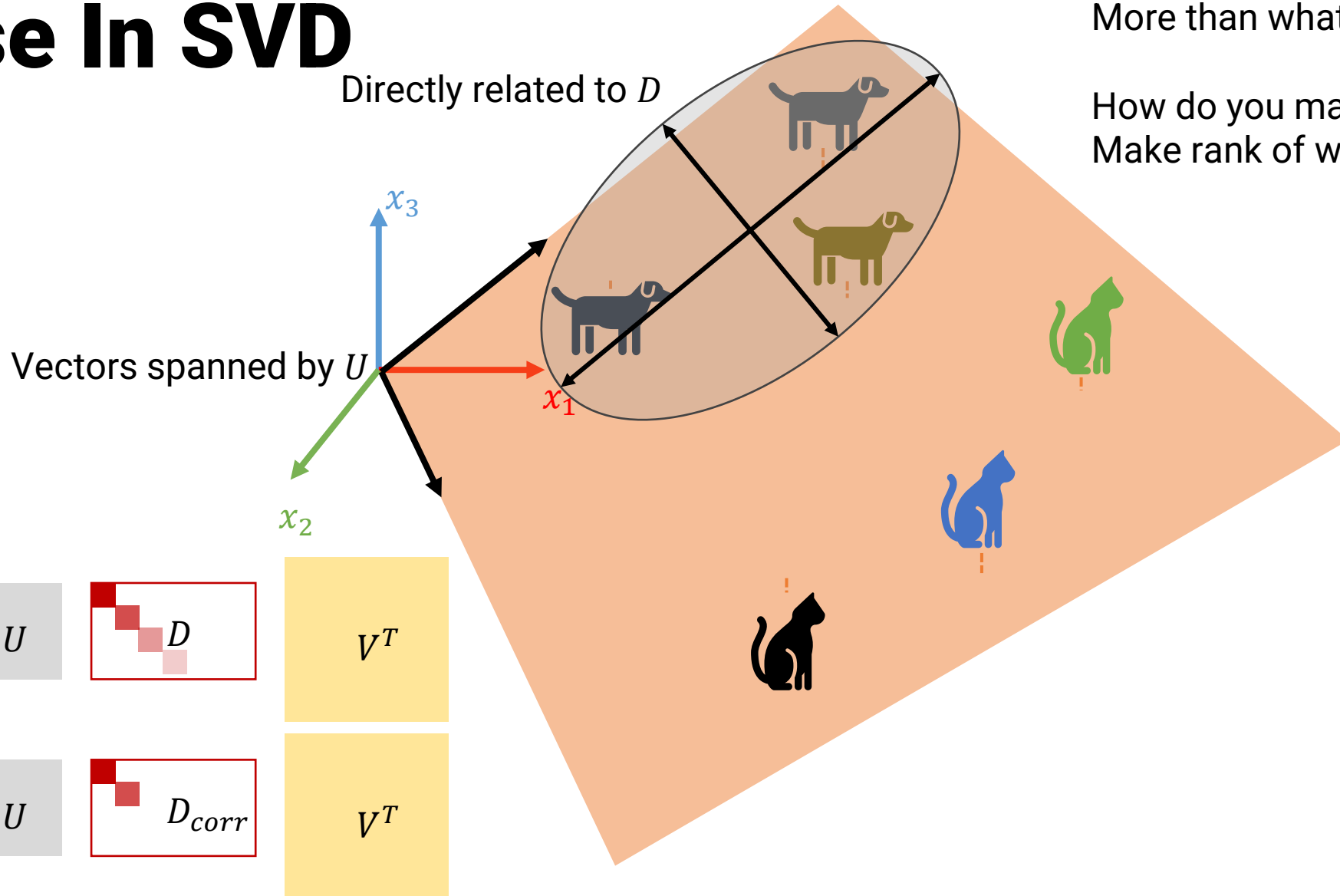




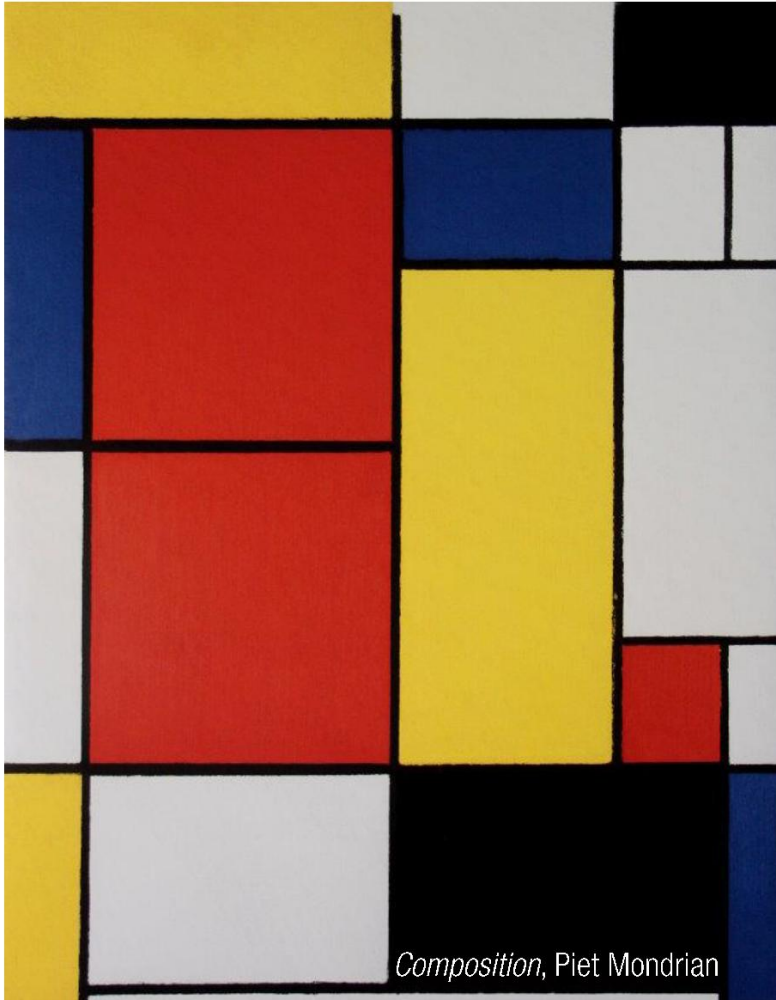
# Noise In SVD

Not all one plane!  
What will be rank of  $D$ ?  
More than what it should be!

How do you make it work?  
Make rank of what it should be!

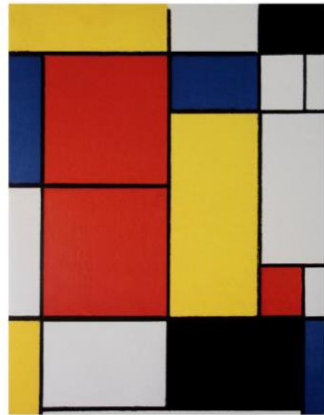


# SVD On Images

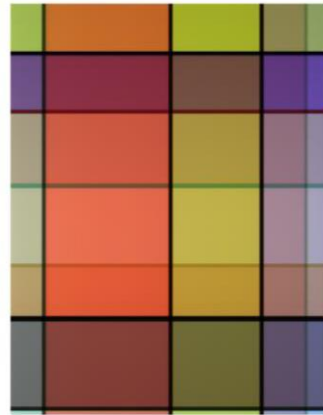


$$= \begin{matrix} \text{U} \\ m \times m \end{matrix} \begin{matrix} \text{D} \\ m \times n \end{matrix} \begin{matrix} \text{V}^T \\ n \times n \end{matrix}$$

# SVD On Images

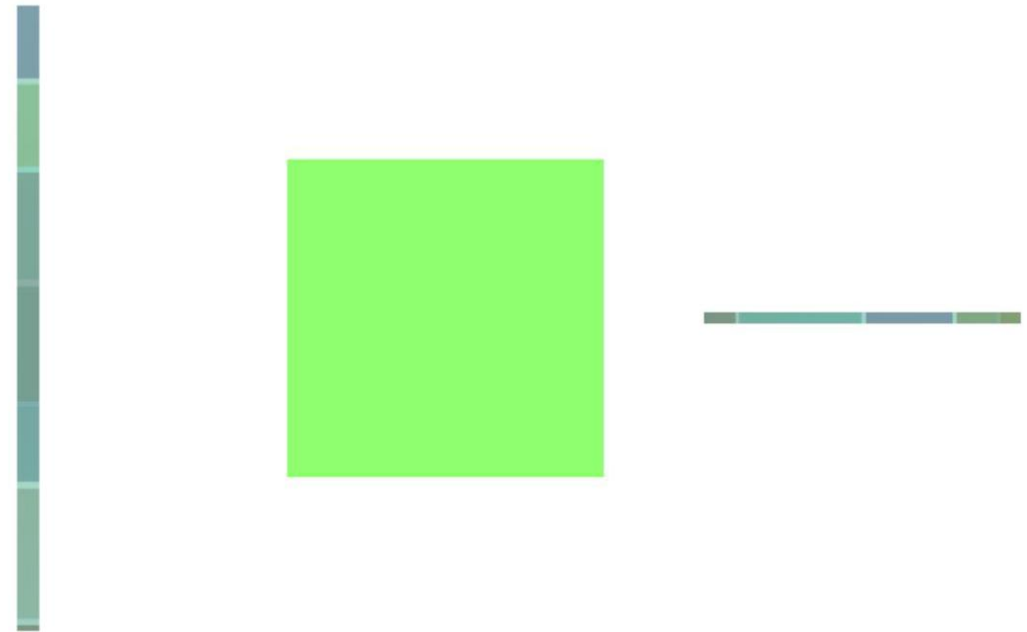


Ground truth



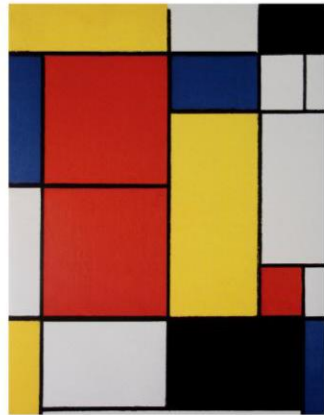
Number of basis: 1

MondrianSVD.m

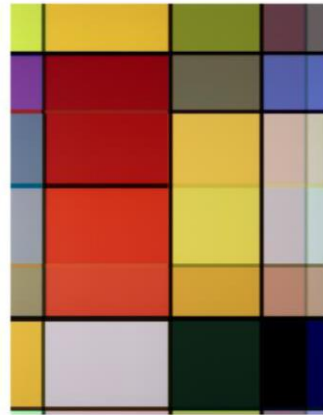


$$A = U D V^T$$

# SVD On Images

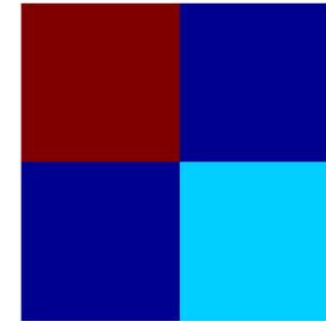


Ground truth



Number of basis: 2

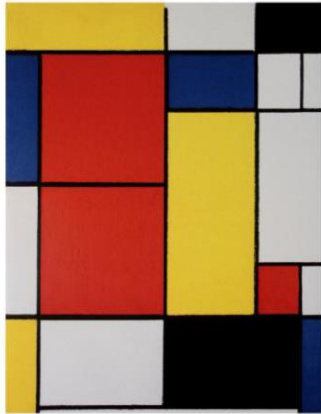
MondrianSVD.m



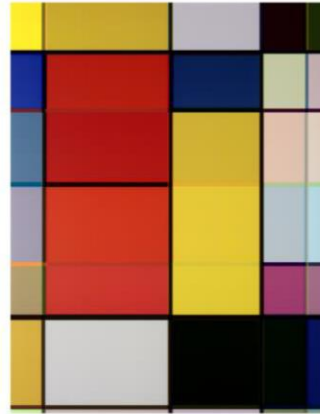
$$A = U D V^T$$

# SVD On Images

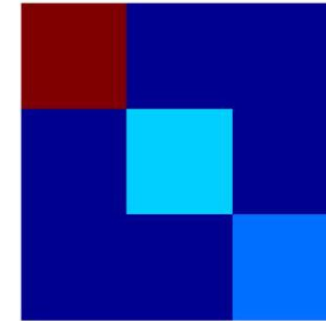
MondrianSVD.m



Ground truth



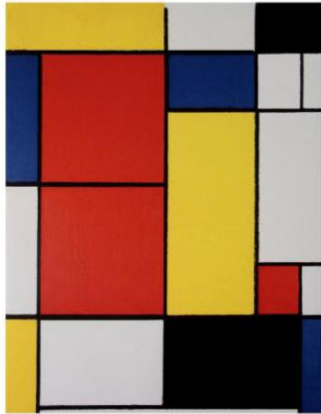
Number of basis: 3



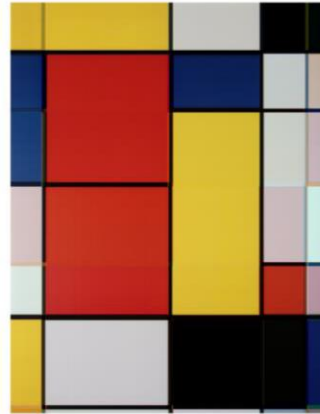
$$A = U D V^T$$

# SVD On Images

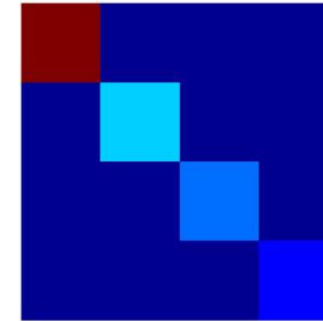
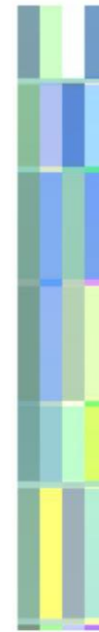
MondrianSVD.m



Ground truth



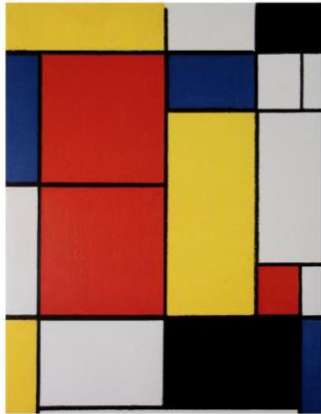
Number of basis: 4



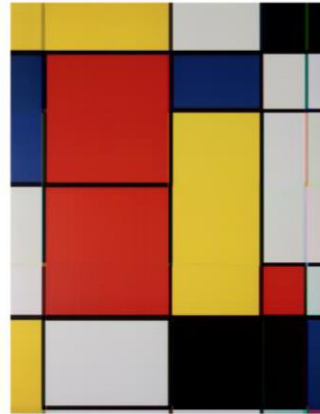
$$A = U D V^T$$

# SVD On Images

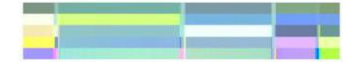
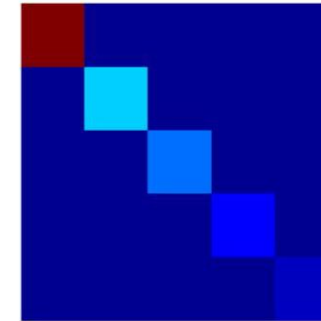
MondrianSVD.m



Ground truth



Number of basis: 5

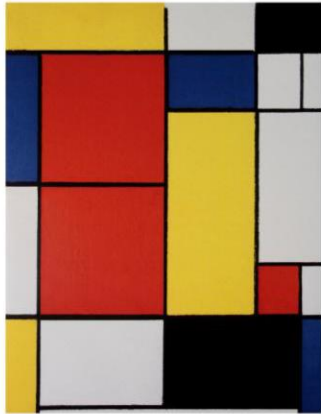


$$A = U D V^T$$

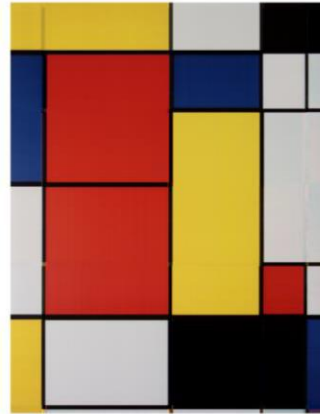


# SVD On Images

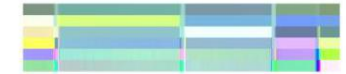
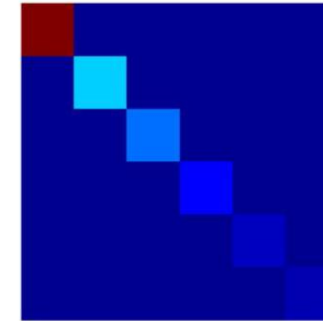
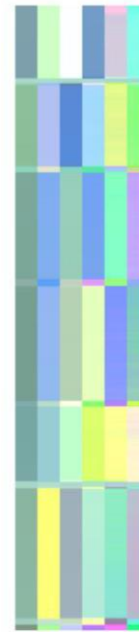
MondrianSVD.m



Ground truth



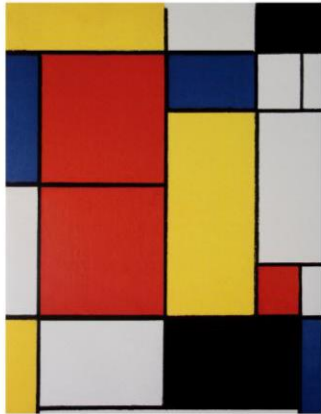
Number of basis: 6



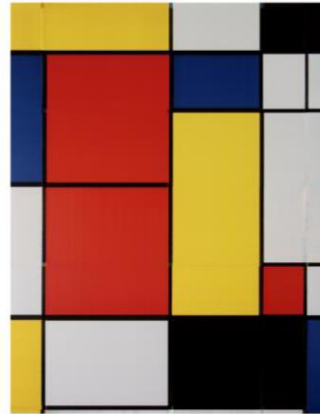
$$A = U D V^T$$

# SVD On Images

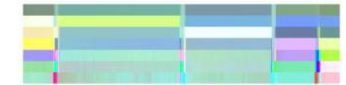
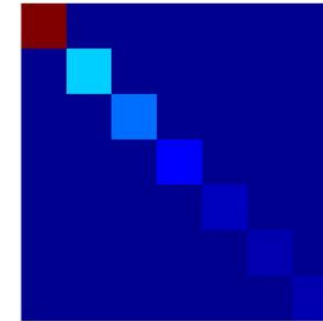
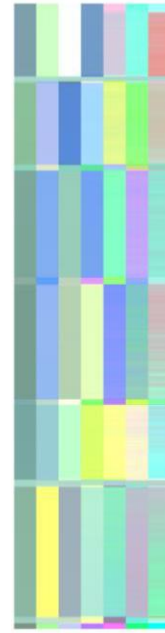
MondrianSVD.m



Ground truth

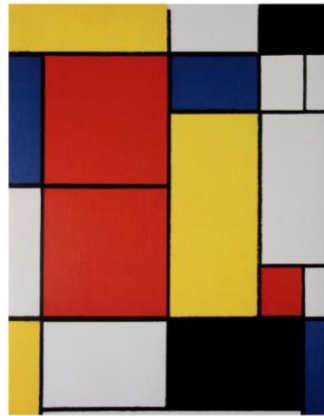


Number of basis: 7



$$A = U D V^T$$

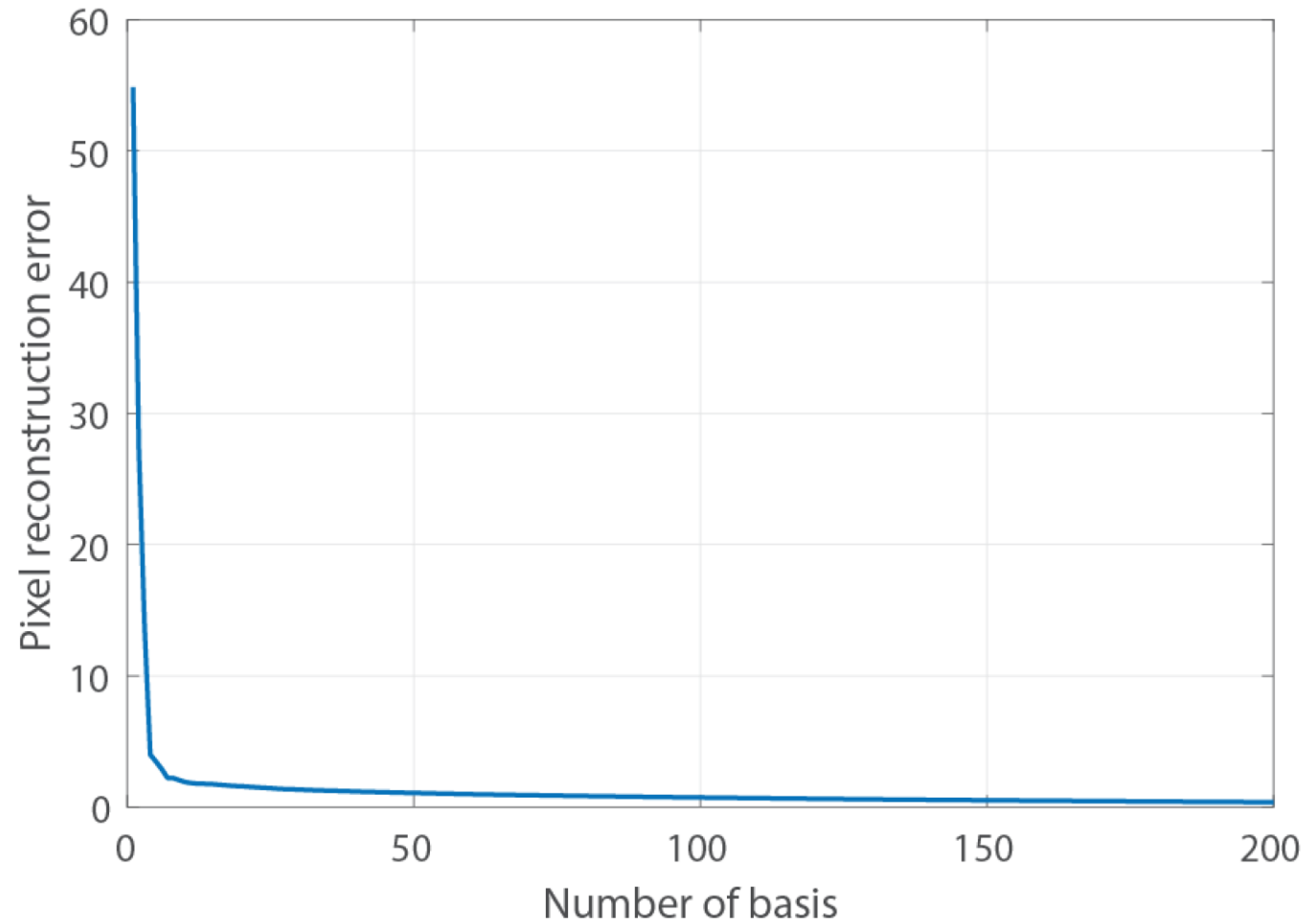
# Reconstruction Error



Ground truth



Number of basis: 7



A

How do you get  $p$  from  $Y$ ?  
You need some function!  
 $p = g(\tilde{Y})$

# Back To Linear Regression

$$Y = W^T \begin{bmatrix} X \\ 1 \end{bmatrix}$$

$$Y = AX$$

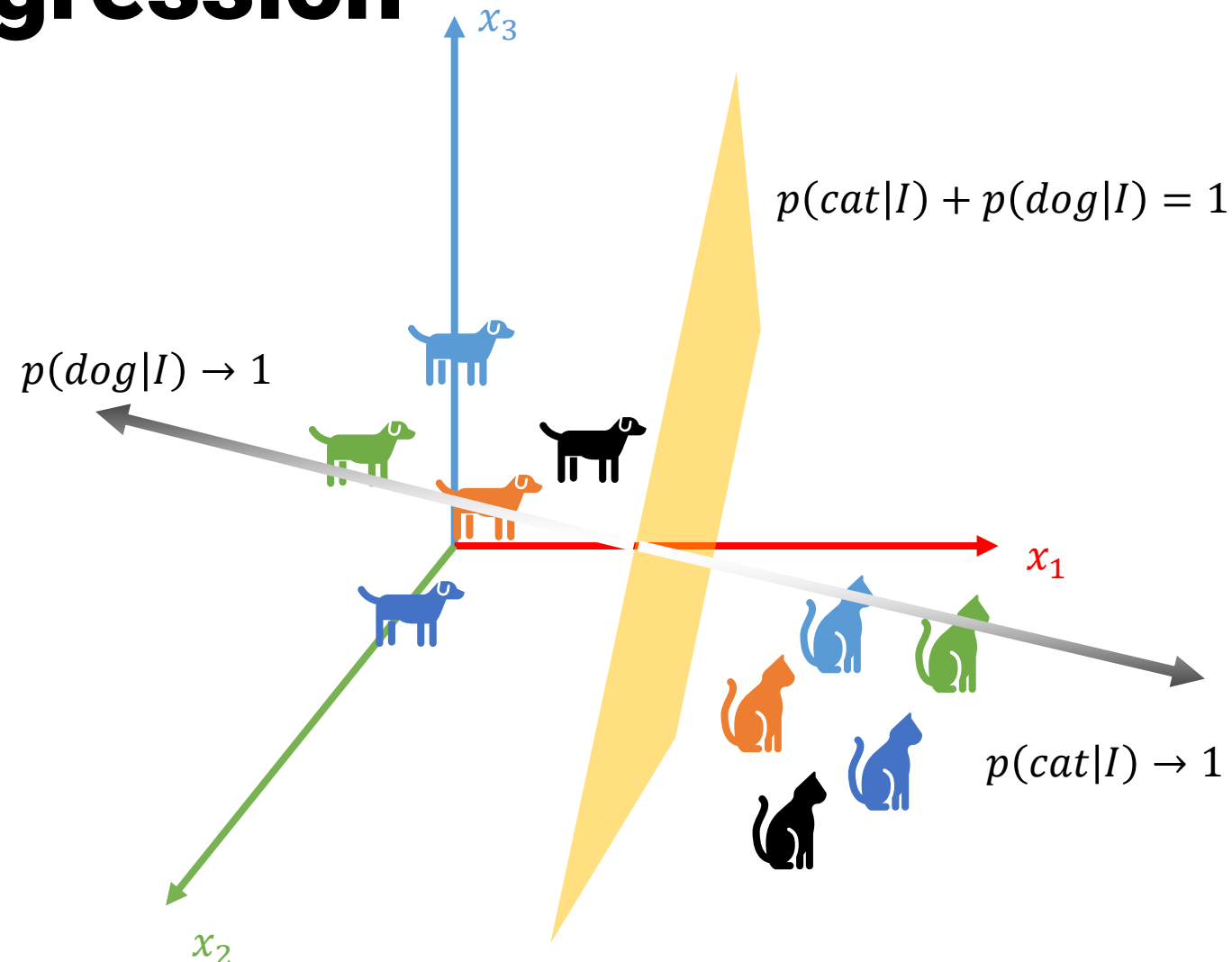
$$X = A^{-1}Y$$

When is this valid?  
 $A$  is square and invertible!

What if  $A$  is not square?  
 $A^{-1} \approx A^+$

How do you get the pseudo-inverse?  
SVD!

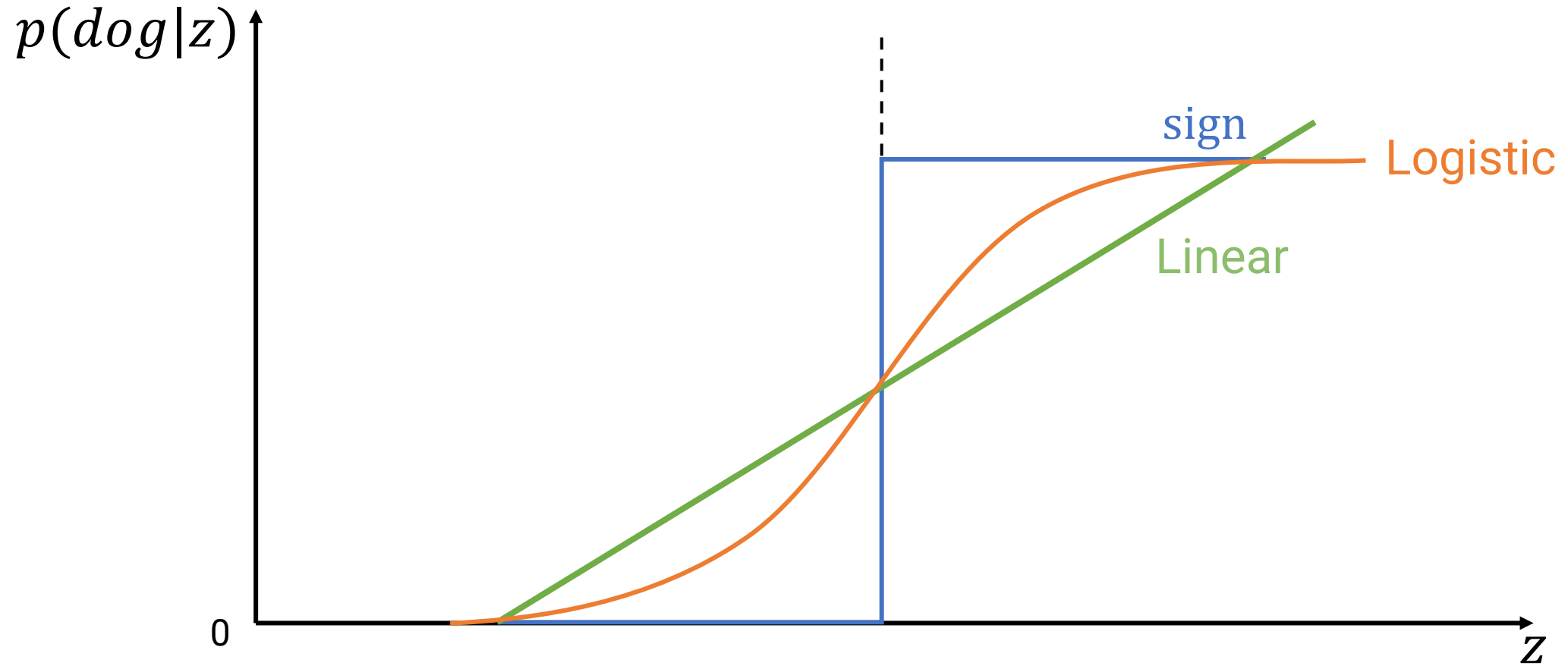
$$A^+ = VD^+U^*$$



$$\tilde{Y} \rightarrow p$$
$$\mathbb{R} \rightarrow [0,1]$$

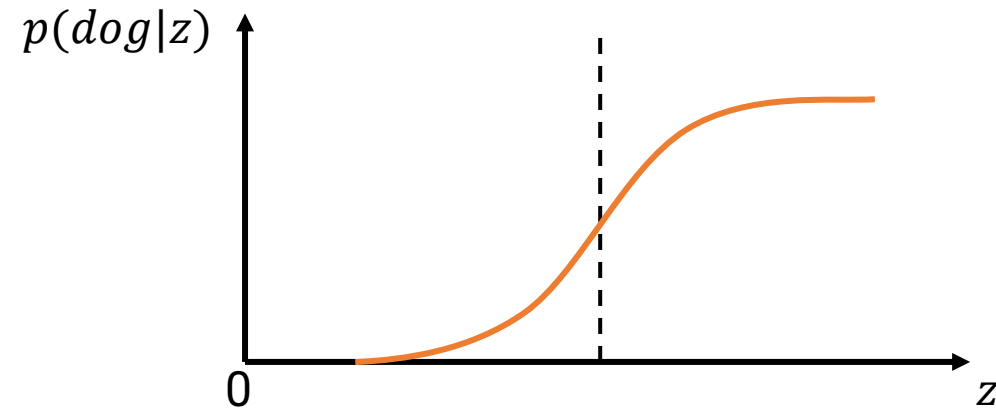
# Function To Get Probabilities

Some squashing function

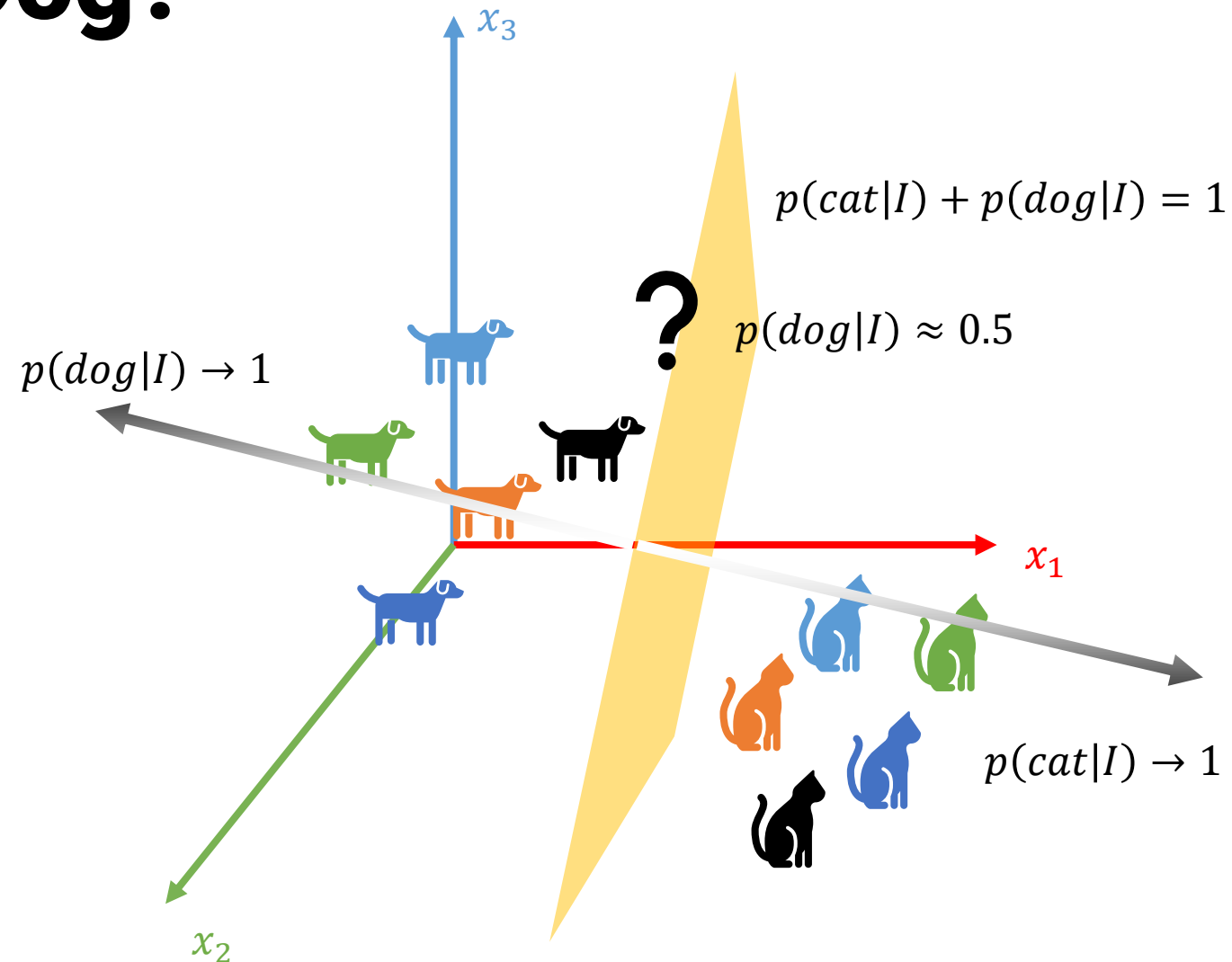


# Logistic Regression

- **Sigmoid function** to obtain probabilities of classification
- Mathematically,  $p = \frac{1}{1+e^{-(W^T X)}}$
- Both class probabilities are related using Log-Odds  
 $z = \log \frac{p}{1-p} = W^T X$  (Also called **Logits**)
- Sometimes also called **Perceptron**

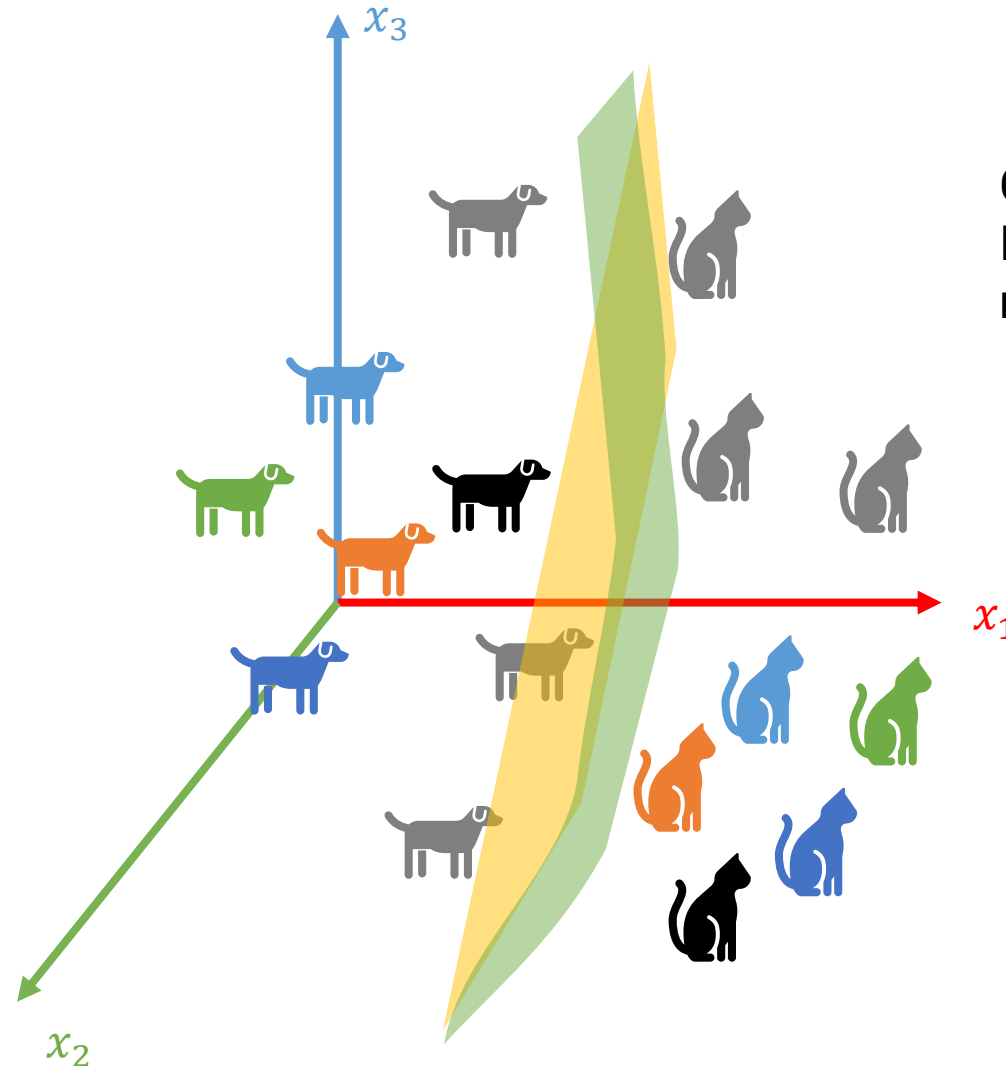


# Is This a Cat or a Dog?



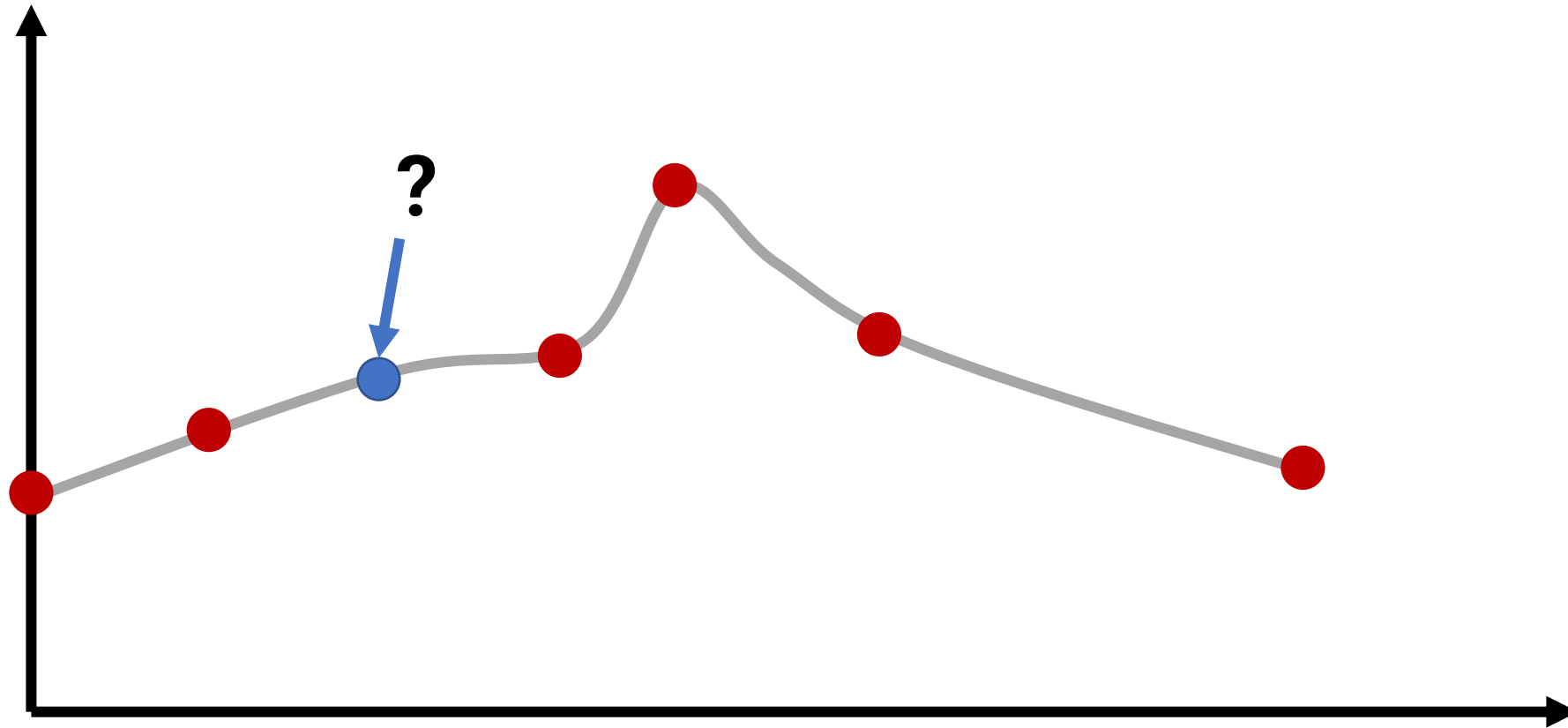


# Why Did This Happen?



Classification boundary is **not linear**!  
It's actually a non-linear function  
map or a **manifold**!

# Interpolation



# Radial Basis Function (RBF)

Instead use a Kernel



From Machine Learning/Math lingo!

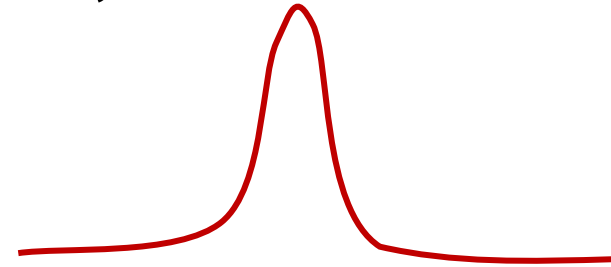
- Smooth
- “Simple”
- Nice properties you want!

Mathematically,

$$\phi: [0, \infty) \rightarrow \mathbb{R}$$

Only depends on **distance** between input and fixed point!

$$\phi = \hat{\phi}(\|\mathbf{x} - \mathbf{c}\|)$$



What is a distance metric?  
Say “how far things are”!

Mathematically,

$$D(\mathbf{x}, \mathbf{x}) = 0$$

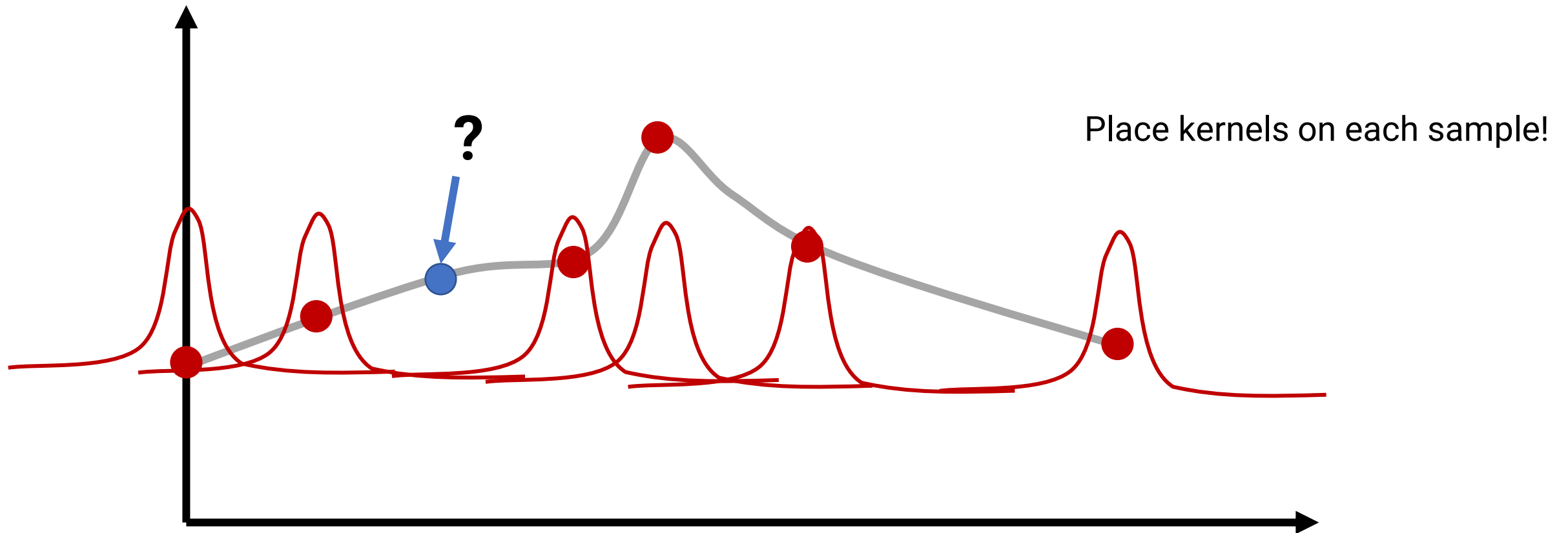
$$D(\mathbf{x}, \mathbf{y}) \geq 0$$

$$D(\mathbf{x}, \mathbf{y}) = D(\mathbf{y}, \mathbf{x})$$

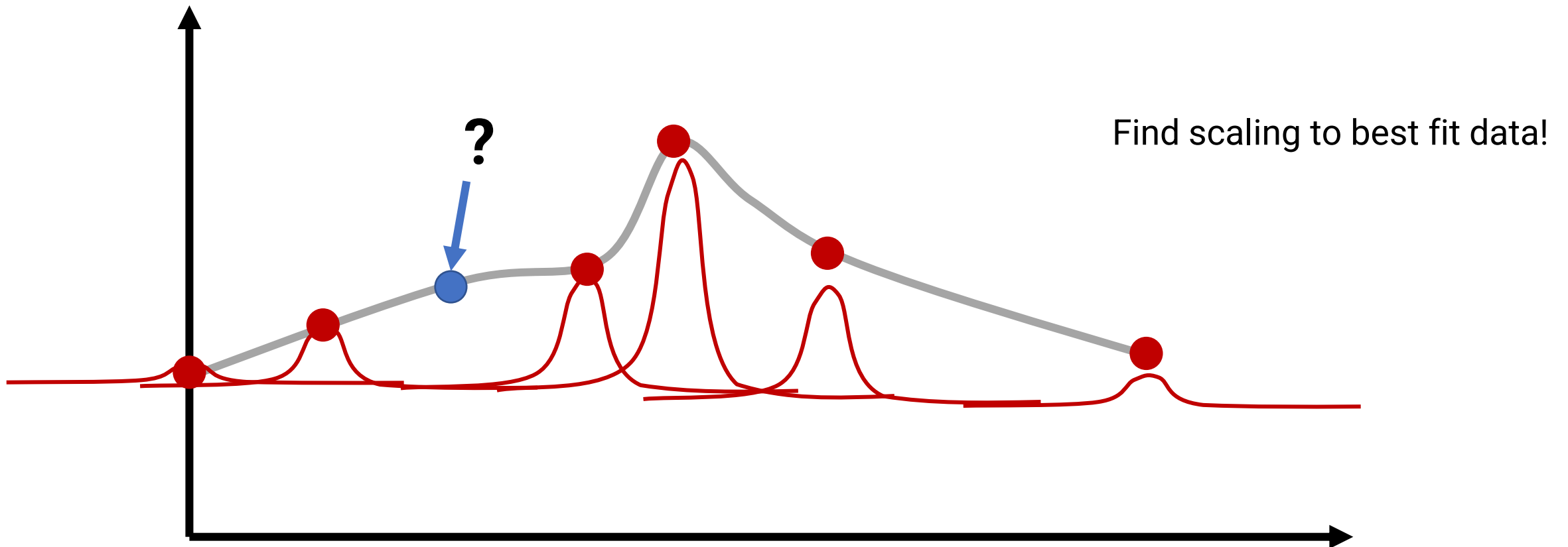
$$D(\mathbf{x}, \mathbf{z}) \leq D(\mathbf{x}, \mathbf{y}) + D(\mathbf{y}, \mathbf{z}) \text{ (Triangle Inequality)}$$

Common distance metrics?  
Norms!

# How To Interpolate?



# How To Interpolate?

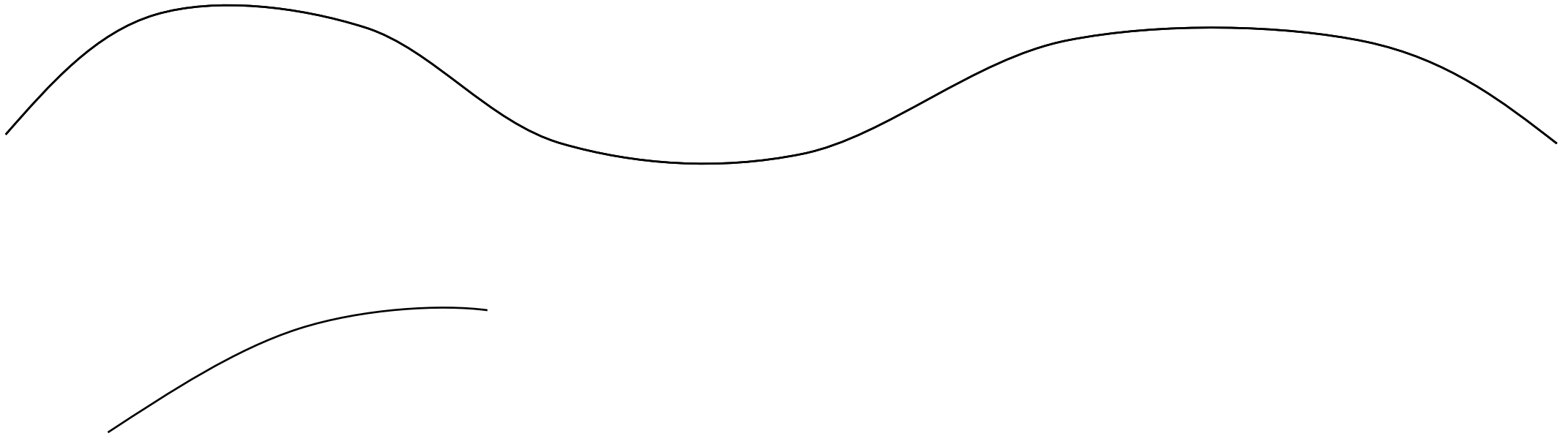


$$f(z) = \sum_{\forall i} \alpha_i R(z, x_i) \text{ given } f(x_i) = y_i$$

# Take A Step Back!



# Take A Step Back!

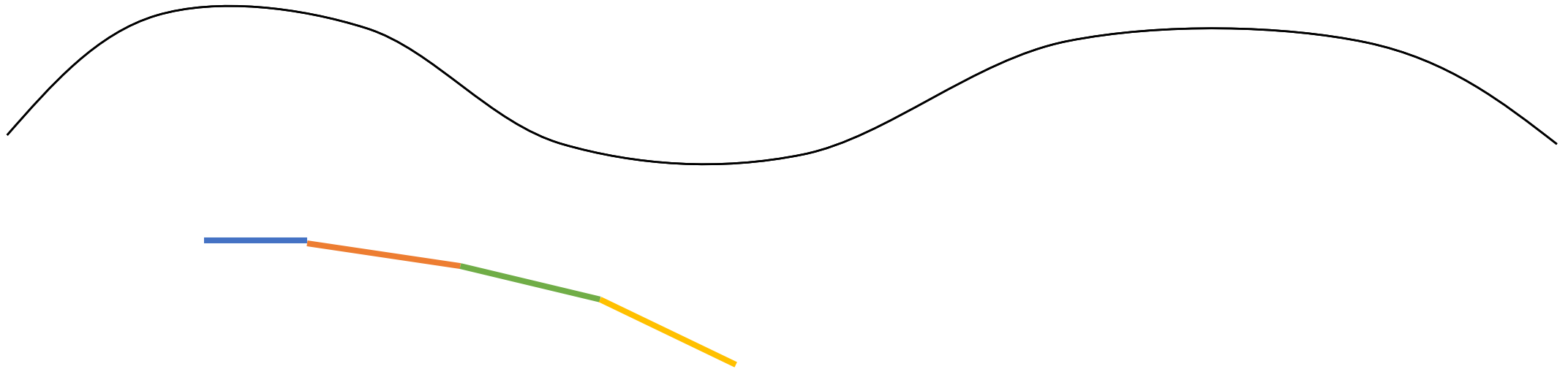


# Take A Step Back!





# Take A Step Back!



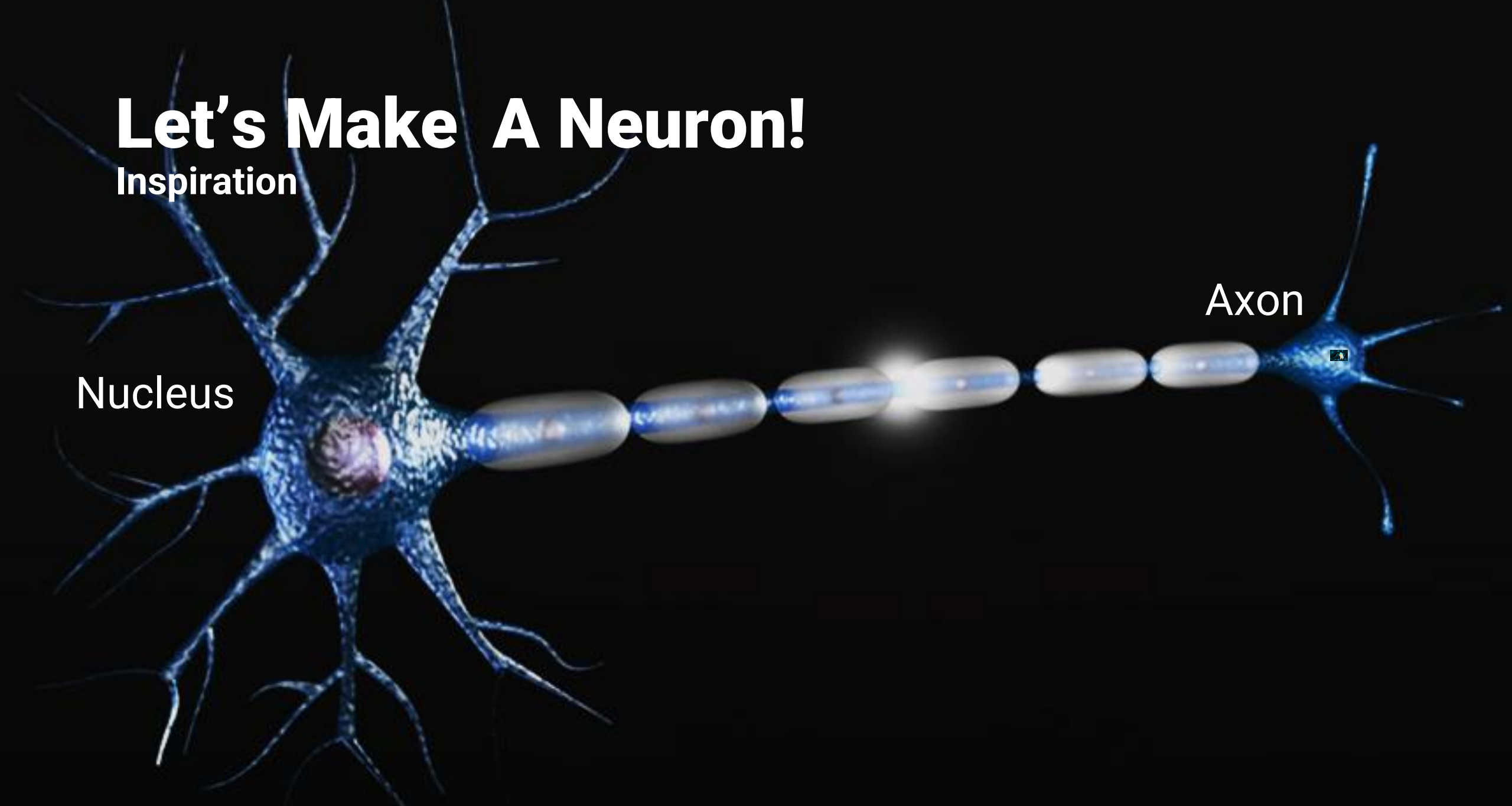
Can approximate any function with multiple simpler functions!

# Let's Make A Neuron!

Inspiration

Nucleus

Axon

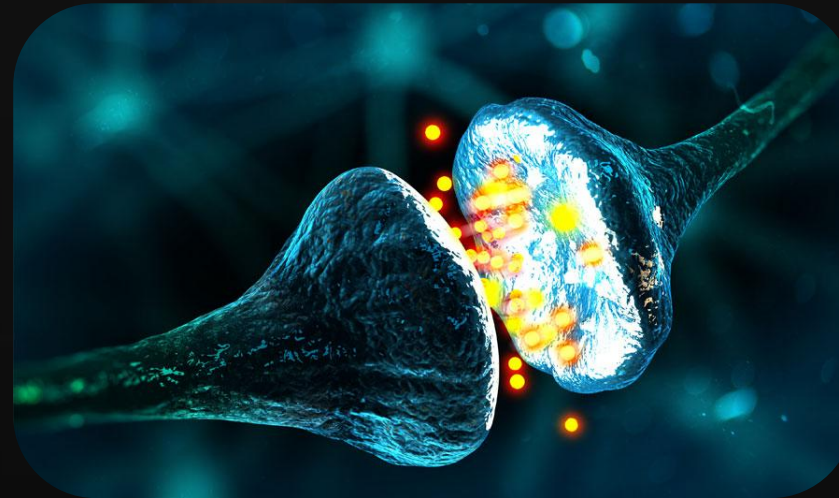


# Let's Make A Neuron!

Inspiration

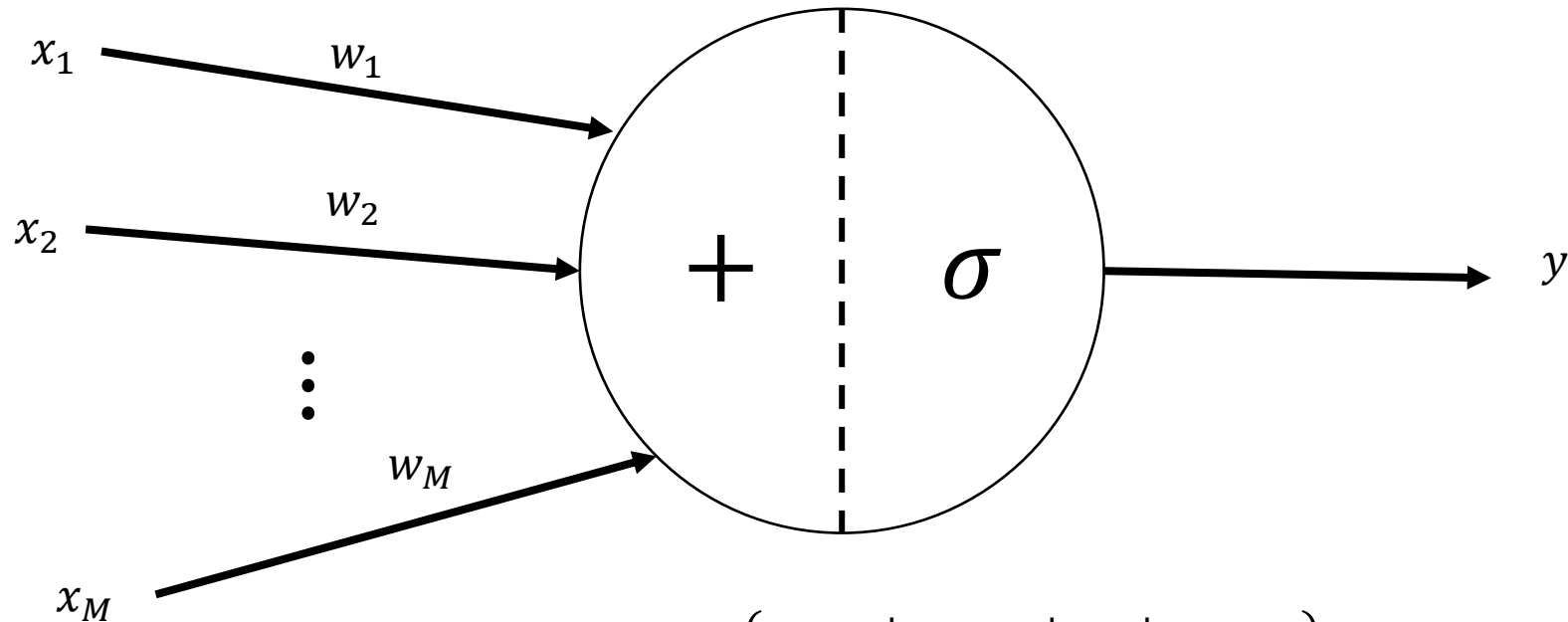
Nucleus

Axon



# Let's Make A Neuron!

Artificial Neuron AKA Perceptron



$$y = \sigma(w_1x_1 + w_2x_2 + \cdots + w_Mx_M)$$

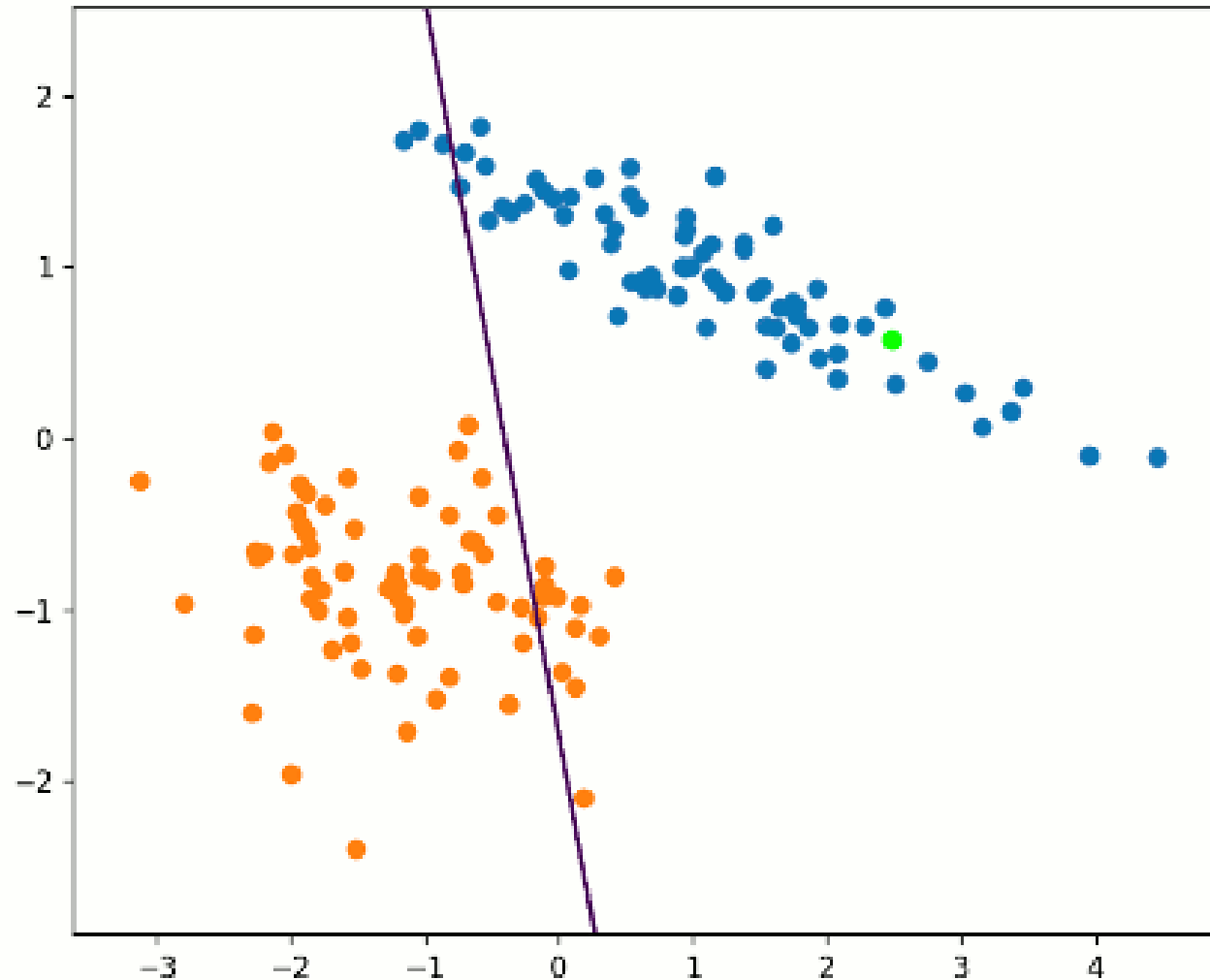
$$y = \sigma(W^T X) \quad \text{Recall Logistic Regression!}$$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_M \end{bmatrix} \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_M \end{bmatrix}$$

# Perceptron Update

McCulloch–Pitts Neuron

Train

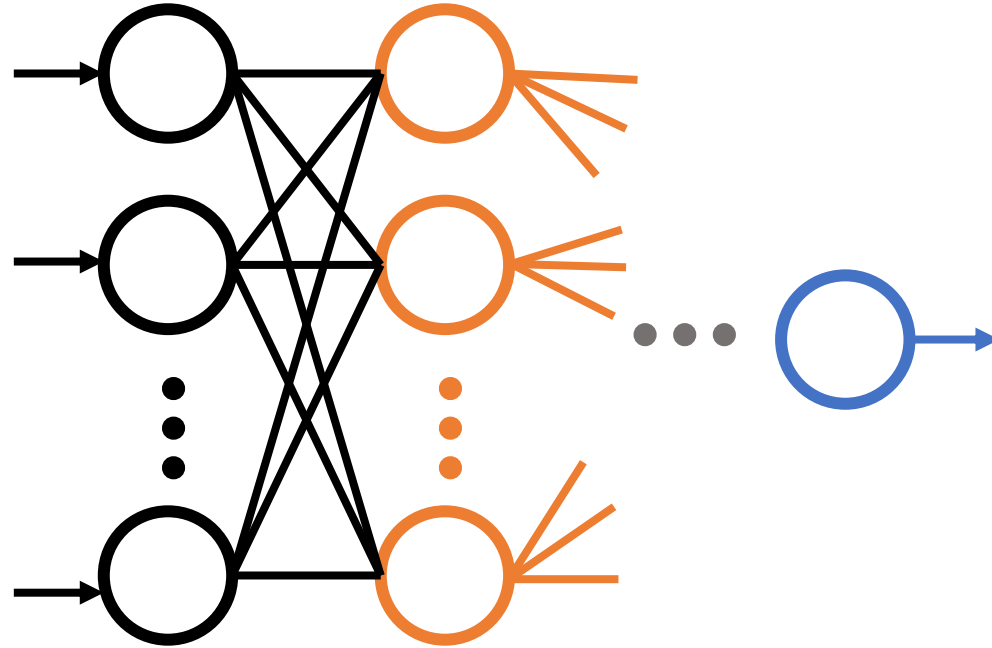


For each train sample  $x_i$

- Compute  $y_i = \sigma(W^T X_i)$
  - Update weight  $W(t + 1) = W(t) + \alpha(\hat{y}_i - y_i)x_i$
- Learning Rate
- Desired Output/Label



# Let's Make A Neural Network!



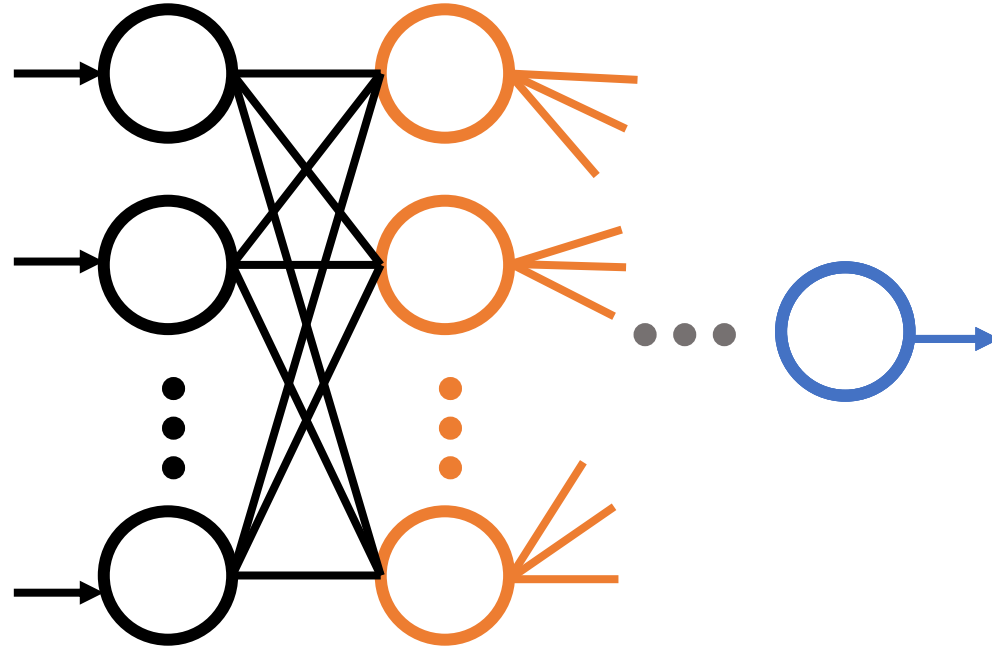
# If $\sigma$ was linear!

$$\begin{aligned}y_1 &= \sigma(W_1 x_1) = A W_1 x_1 \\y_2 &= \sigma(W_2 y_1) = \sigma(\sigma(W_1 x_1)) = A^2 W_2 W_1 x_1 = W_2' x_1 \\y_n &= \prod_{\forall i} y_{n-i} \cdots A W_1 x_1 = W_n' x_1\end{aligned}$$

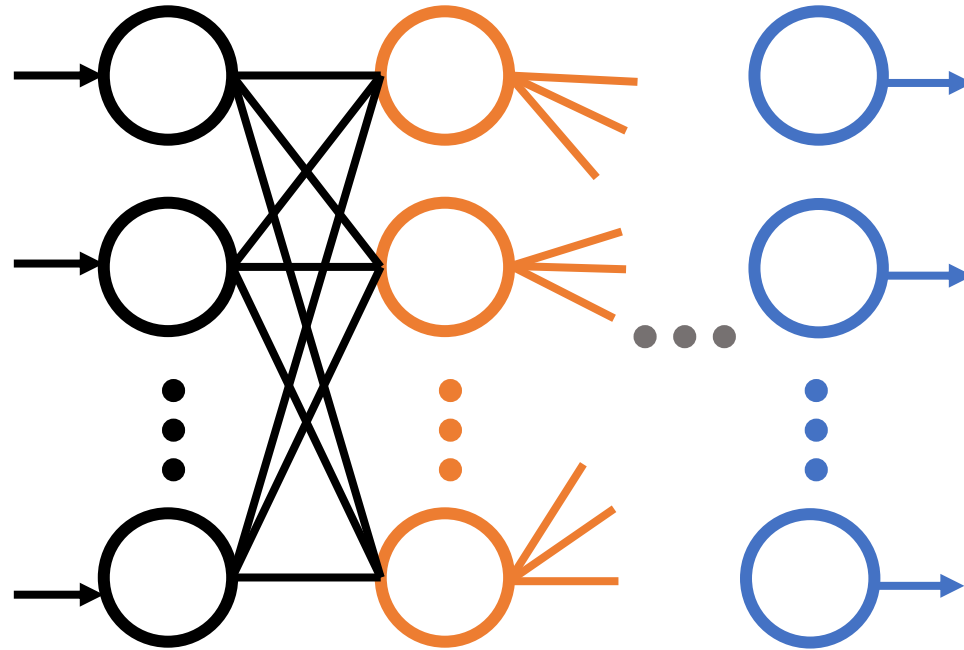
**This is the same as one neuron!**



# Let's Make A Neural Network!



# Let's Make A Neural Network!



# Training-Validation-Test Sets

- **Training:** A labeled group of examples which is fed to the NN during training (weight updating/backprop etc.)
- **Validation:** A group of unseen examples used to test your NN to see how it performs
- **Test:** The set where your model is deployed on!



dog (1)



dog (1)



dog (1)



cat (0)



dog (1)



dog (1)



cat (0)

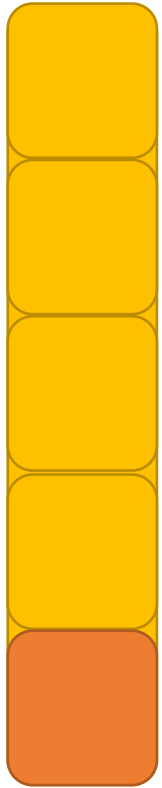


cat (0)



dog (1)

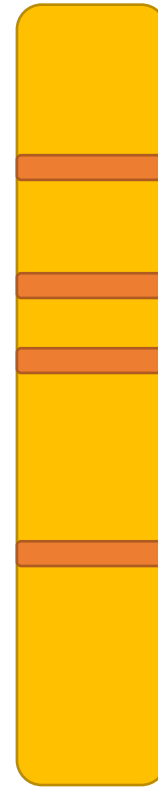
# Cross-Validation



K-fold



Leave-one-out



Random Sampling



Test Set

# Let's Get Some Math Done!

- Derivative is defined by
- $D_x f = \lim_{\|h\| \rightarrow 0} \frac{f(x+h) - f(x)}{\|h\|}$
- $Df = (\partial_{x_1} f \quad \partial_{x_2} f \quad \partial_{x_3} f)$  (assuming variables are column vectors)
- **Gradient** is the **adjoint** of derivative!

$$Df = (\partial_{x_1} f \quad \partial_{x_2} f \quad \partial_{x_3} f) \text{ and } \nabla f = \begin{bmatrix} \partial_{x_1} f \\ \partial_{x_2} f \\ \partial_{x_3} f \end{bmatrix}$$

For matrices:

$$\nabla f = \begin{bmatrix} \partial_{11} f & \partial_{12} f & \partial_{13} f \\ \partial_{21} f & \partial_{22} f & \partial_{23} f \\ \partial_{31} f & \partial_{32} f & \partial_{33} f \end{bmatrix} \text{ where } x = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix}$$



# What Is The Gradient?

## Function

$$f(x) = Ax$$

$$f(x) = \|x\|^2$$

$$f(x) = \frac{1}{2} x^T Ax$$

Symmetric  $A$

## Gradient

$$\nabla f(x) = A^T$$

$$\nabla f(x) = 2x$$

$$\nabla f(x) = Ax$$



Remember:

$$\text{If } h(x) = f(Ax) \Rightarrow \nabla h(x) = A^T f'(Ax)$$

$$\text{And if } h(x) = f(xA) \Rightarrow \nabla h(x) = f'(xA)A^T$$

# Recall Chain Rule!

$$h(x) = f \circ g(x) = f(g(x))$$

For single variables:

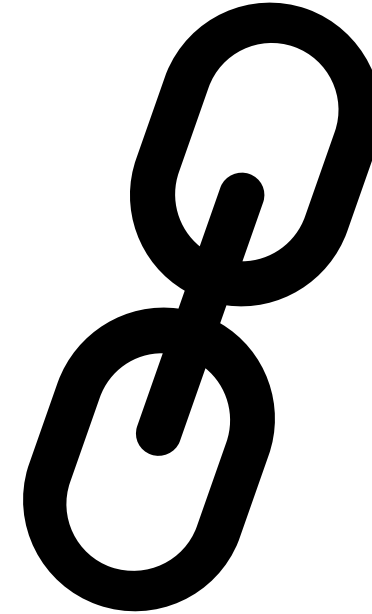
$$\partial f(g(x)) = f'(g)g'(x)$$

For multiple-variables:

$$Df(g(x)) = Df \circ D(g(x))$$

For gradients:

$$\nabla f(g(x)) = \nabla g(x) \nabla f(x)$$



# Recall Ridge Regression

## Least Squares With Penalty

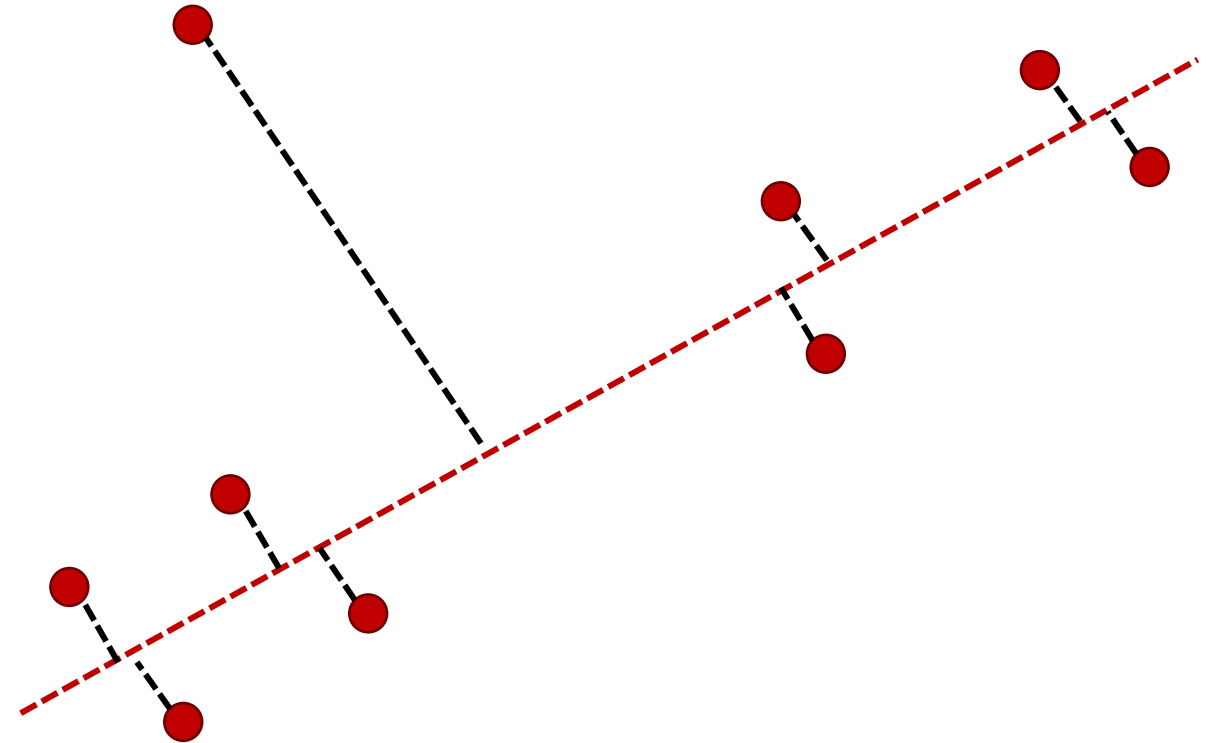
$$f(x) = \frac{1}{2} \|Ax - b\|^2 + \frac{\lambda}{2} \|x\|^2$$

How do you obtain optimal solution?  
Set gradient to zero!

$$\nabla f(x) = \lambda x + A^T(Ax - b) = 0$$

$$(A^T A + \lambda I)x = A^T b$$

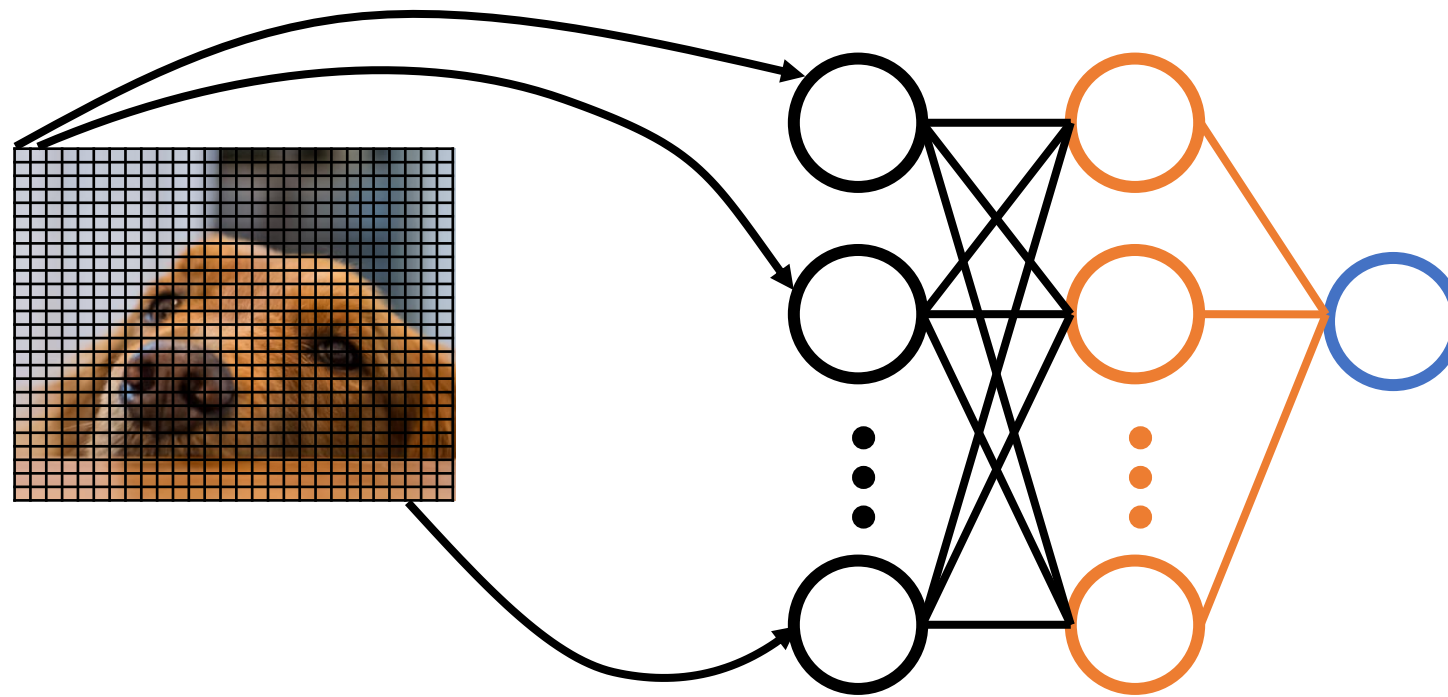
$$x^* = (A^T A + \lambda I)^{-1} A^T b$$





# Key To Training: Error Back Propagation

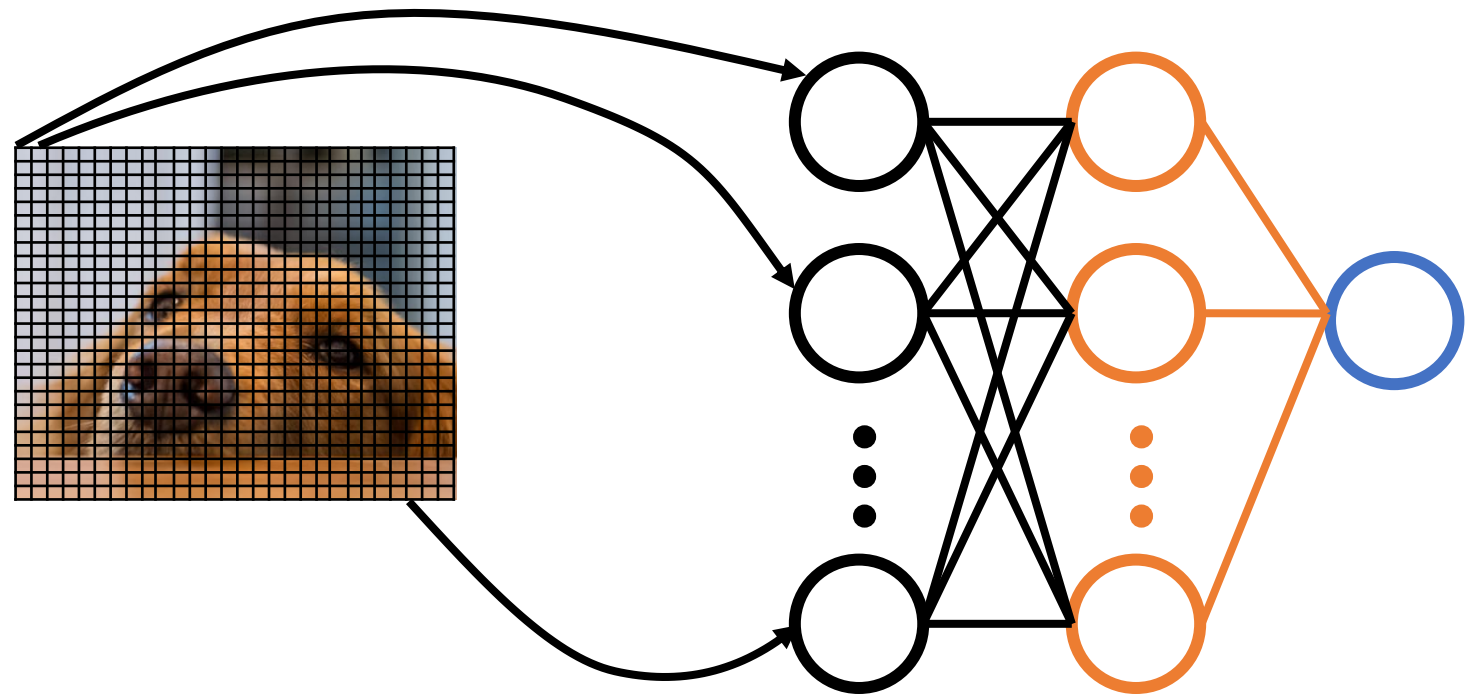
AKA Backprop



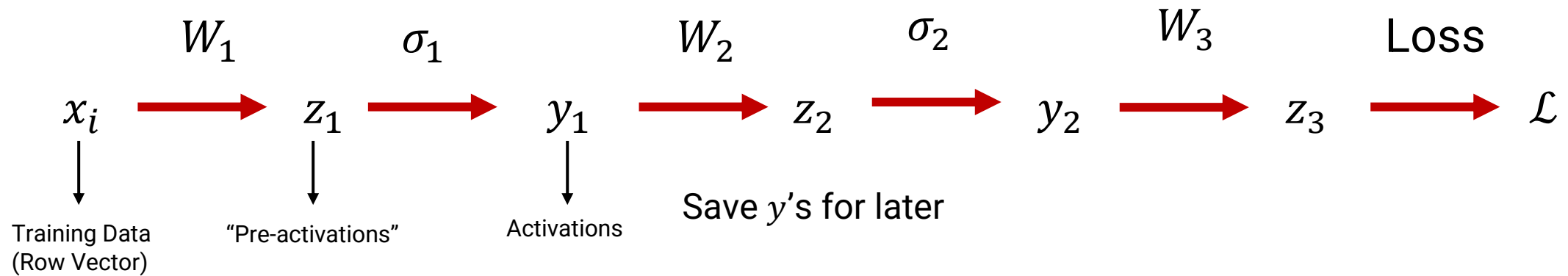
$x_i$   
Input Layer

$y_1 = \sigma(x_i W_1)$   $y_2 = \sigma(y_1 W_2)$   $y_3 = \sigma(y_2 W_3)$   
Hidden Layer 1 Hidden Layer 2 Output Layer

# Forward Pass

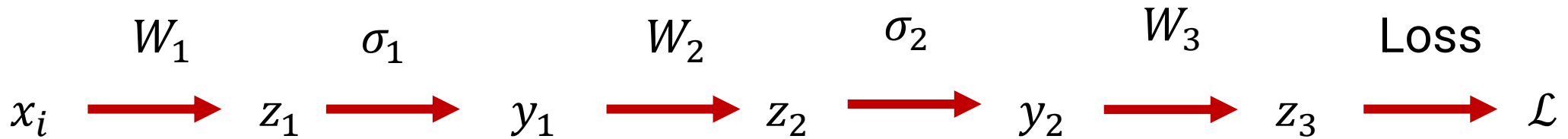
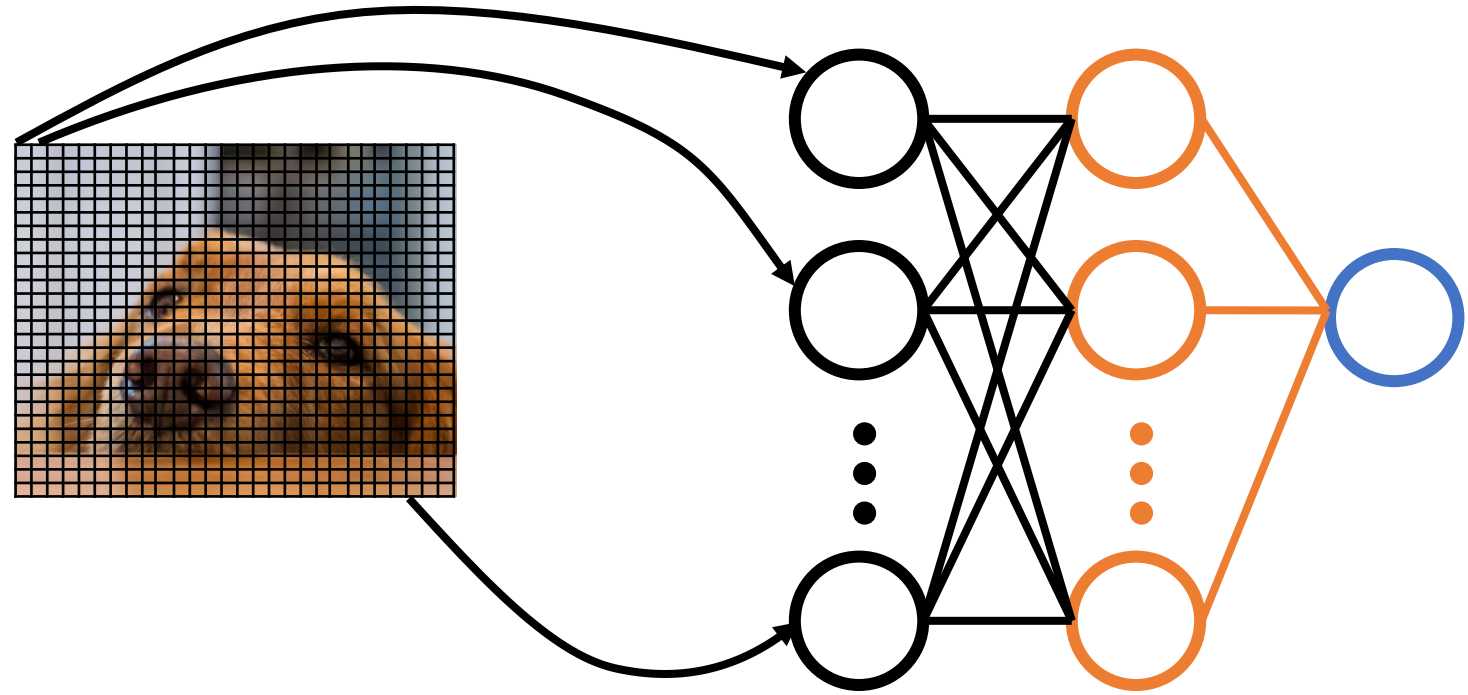


Training Data  
(Row Vector)



# Backward Pass

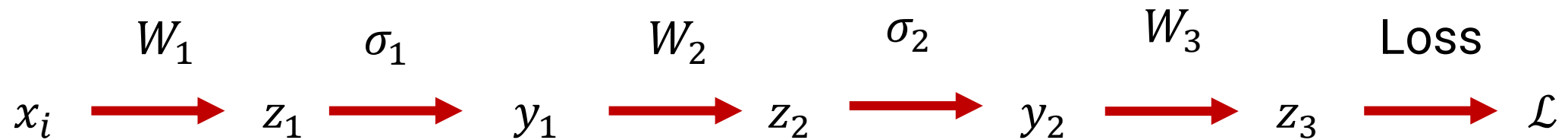
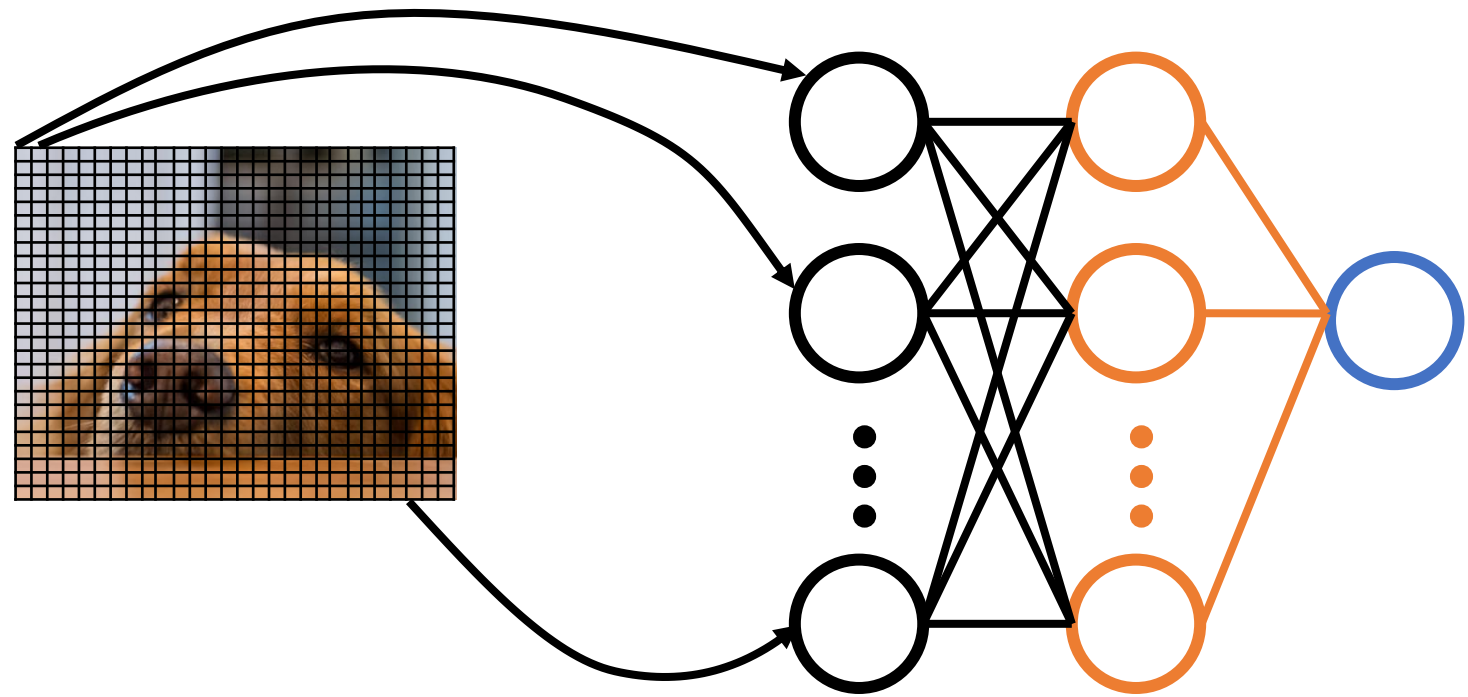
We need derivative of  $\mathcal{L}$  with respect to  $x$  and  $W$



Derivative here is  $\frac{\partial \mathcal{L}}{\partial z_3}$

# Backward Pass

We need derivative of  $\mathcal{L}$  with respect to  $x$

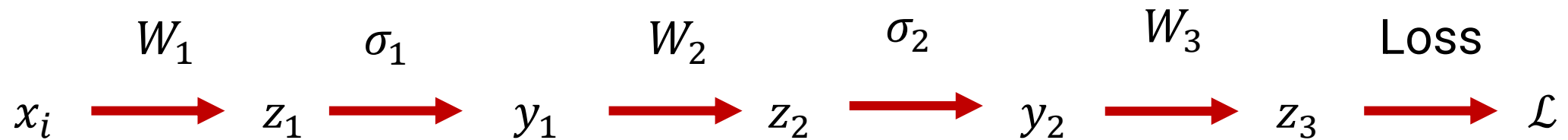
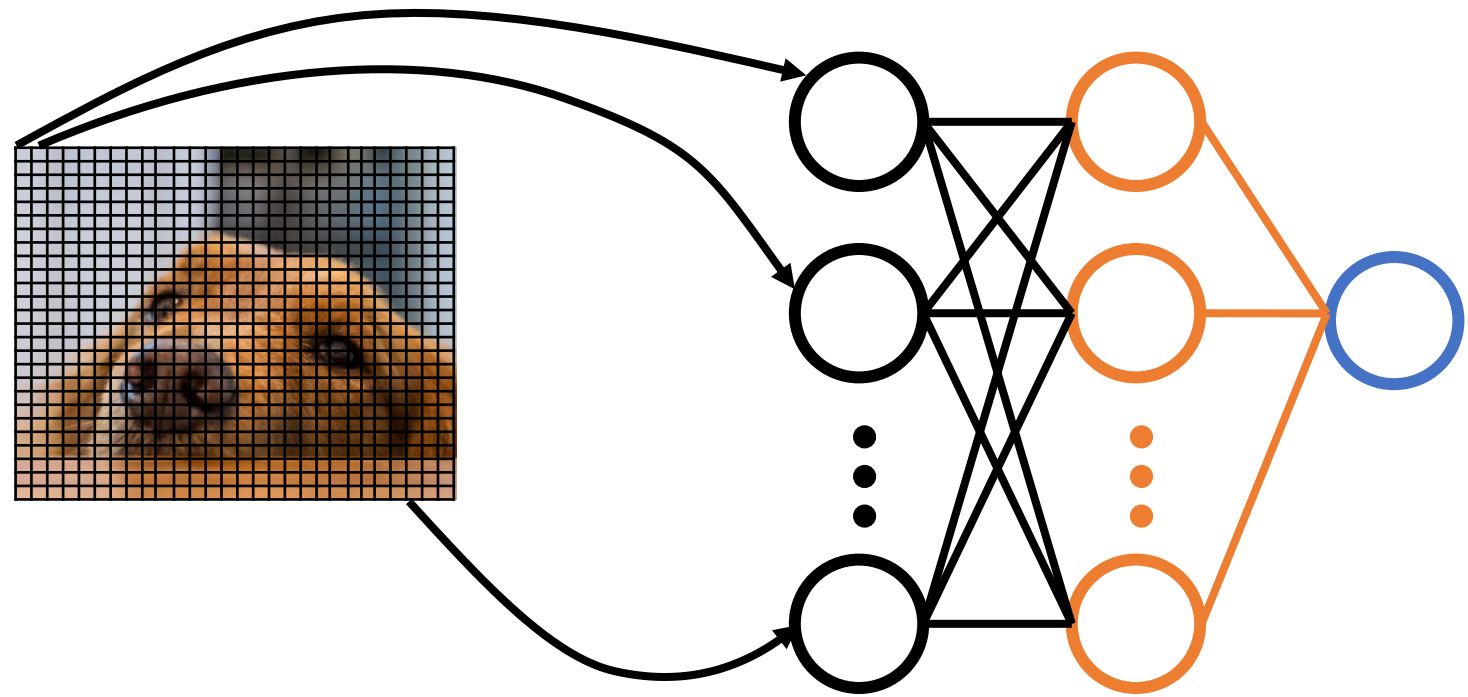


Derivative here is  $W_3 \frac{\partial \mathcal{L}}{\partial z_3}$

Scaling property!

# Backward Pass

We need derivative of  $\mathcal{L}$  with respect to  $x$

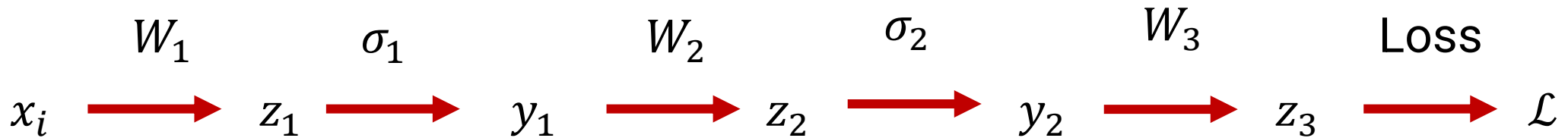
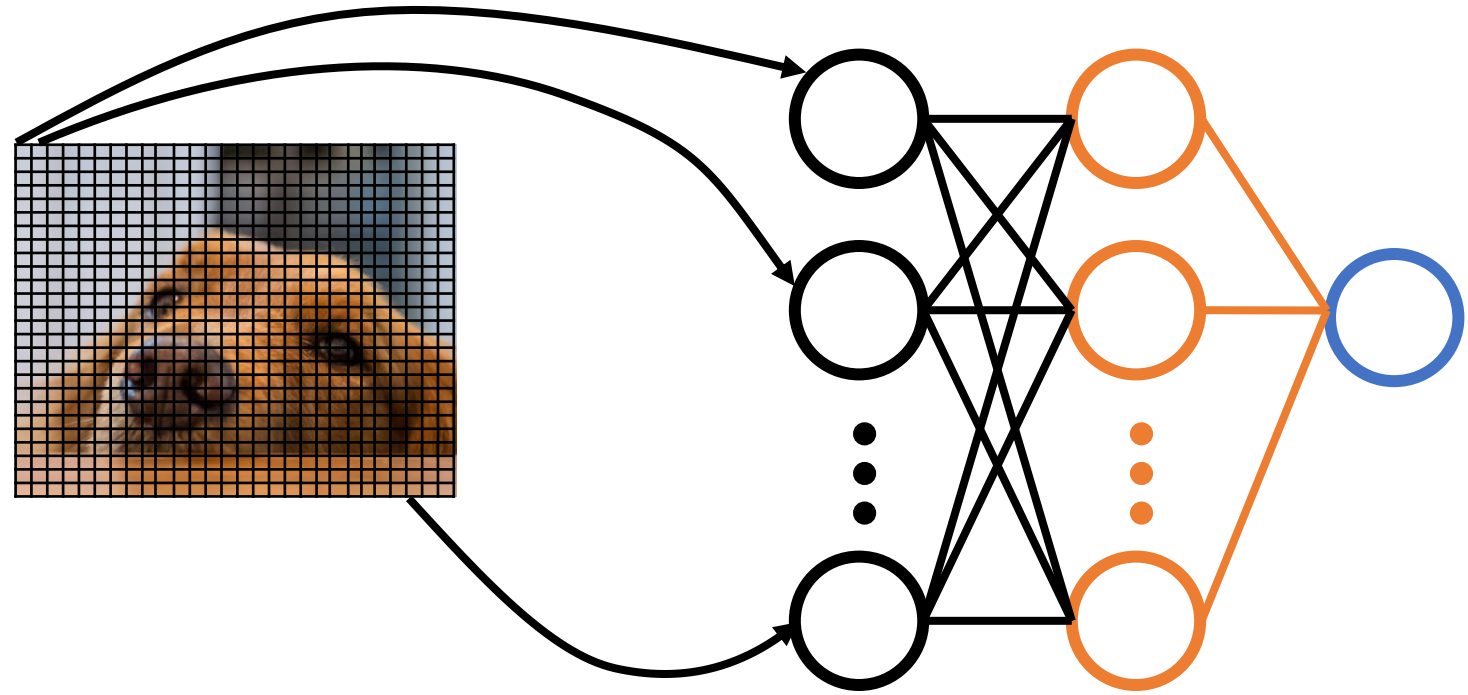


Derivative here is  $\sigma_2' W_3 \frac{\partial \mathcal{L}}{\partial z_3}$

Chain Rule!

# Backward Pass

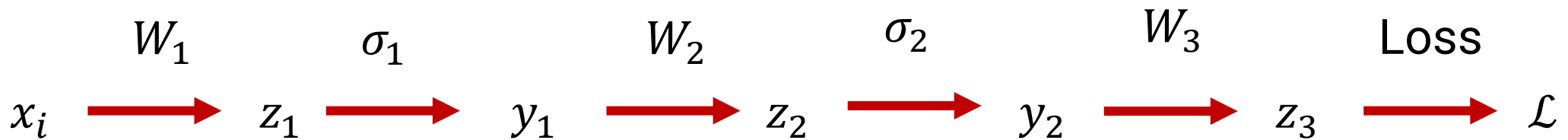
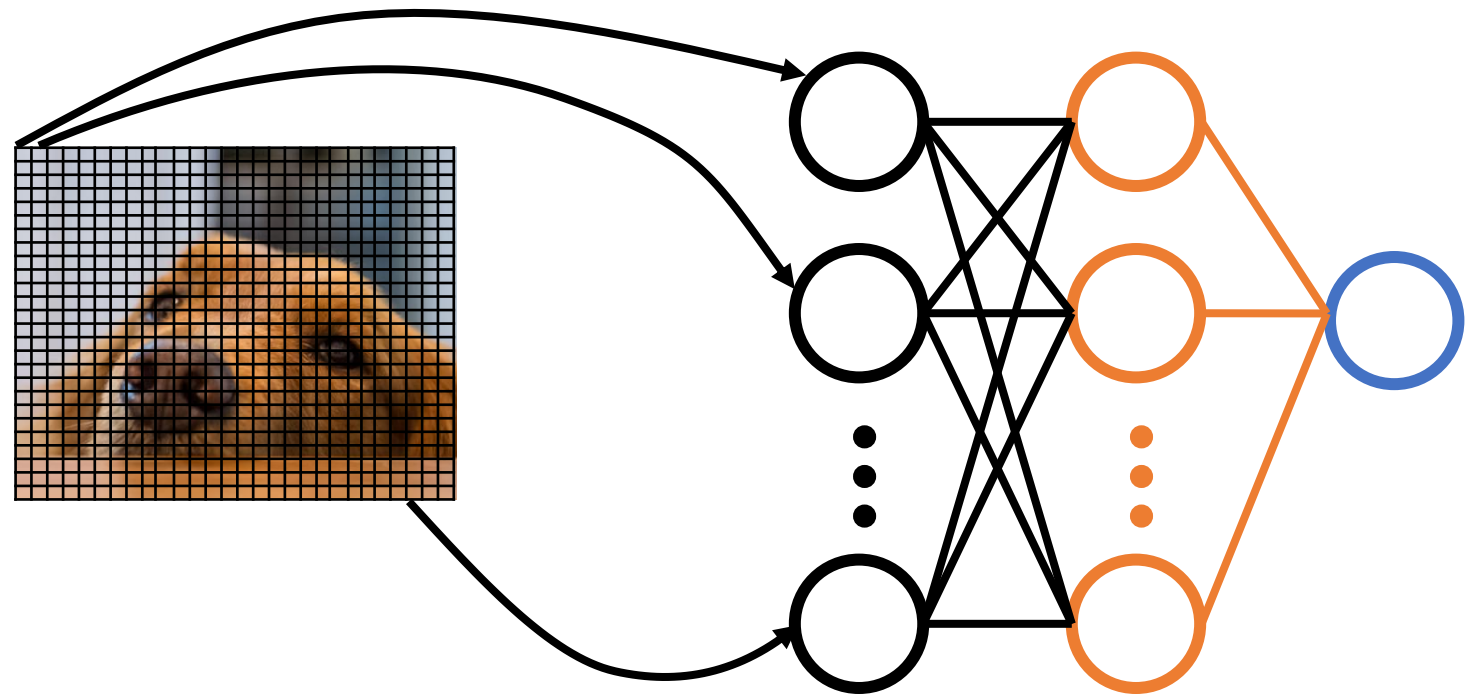
We need derivative of  $\mathcal{L}$  with respect to  $x$



Derivative here is  $W_2 \sigma_2' W_3 \frac{\partial \mathcal{L}}{\partial z_3}$

# Backward Pass

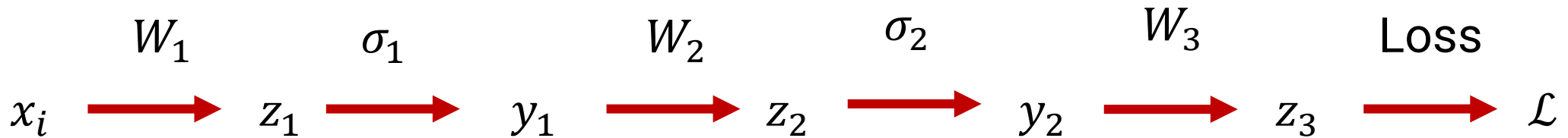
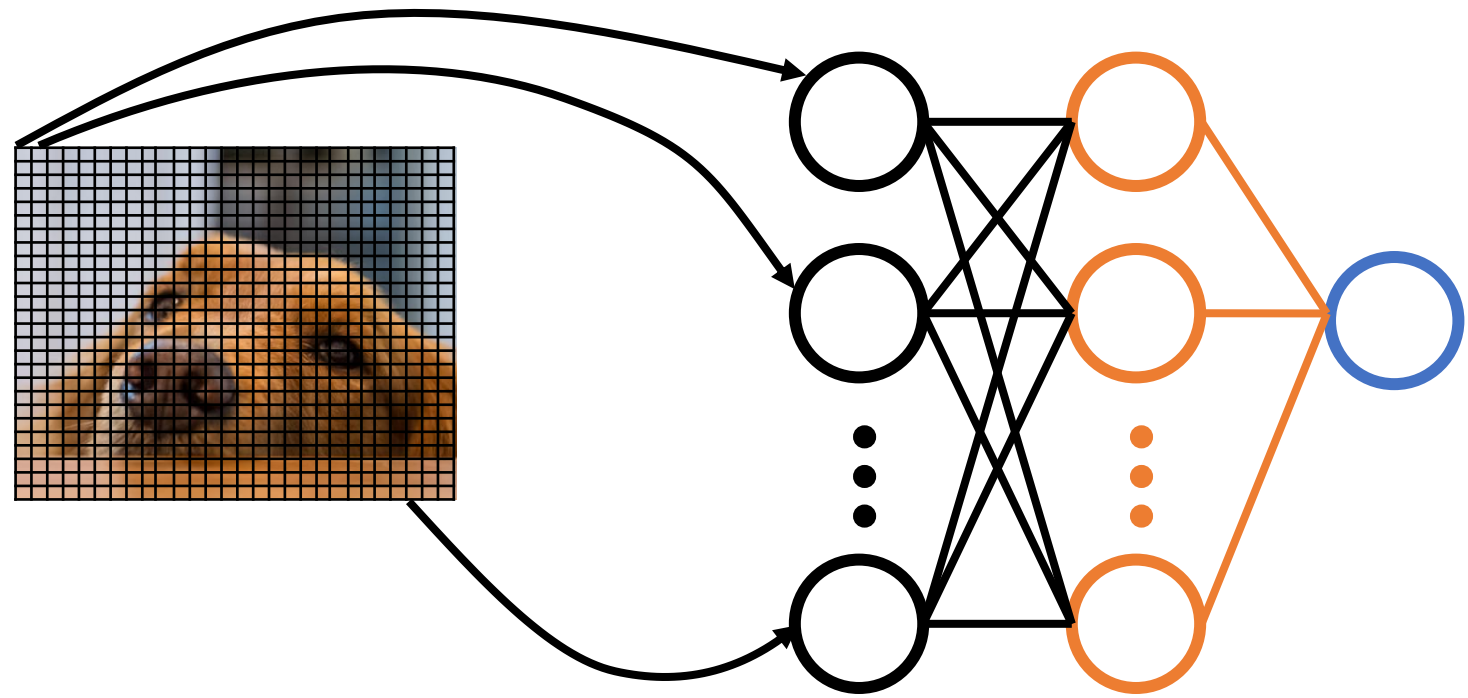
We need derivative of  $\mathcal{L}$  with respect to  $x$



Derivative here is  $\sigma_1' W_2 \sigma_2' W_3 \frac{\partial \mathcal{L}}{\partial z_3}$

# Backward Pass

We need derivative of  $\mathcal{L}$  with respect to  $x$

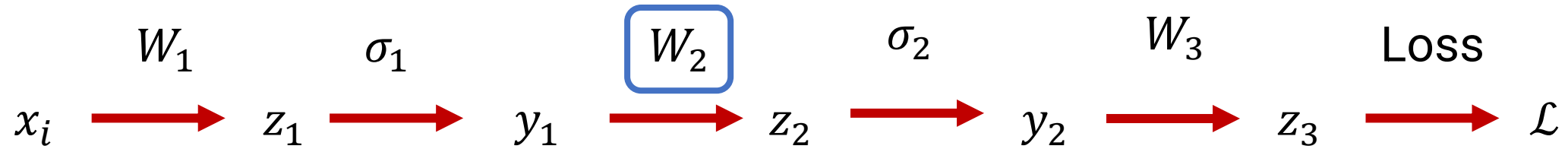


Derivative here is  $W_1 \sigma_1' W_2 \sigma_2' W_3 \frac{\partial \mathcal{L}}{\partial z_3}$



# Derivative With Respect To $W_2$

Derivative of loss  $\mathcal{L}$  with respect to  $W_2$



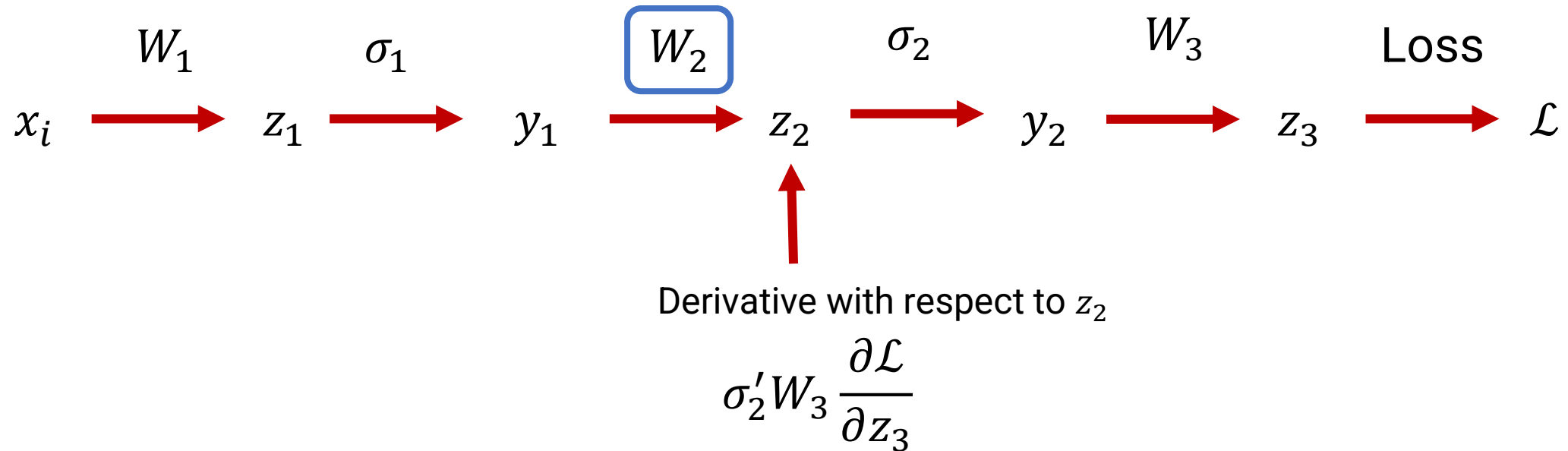
Use Chain Rule

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial W_2}$$

$$\nabla_{W_2} \mathcal{L} = \nabla_{W_2} z_2 \nabla_{z_2} \mathcal{L}$$

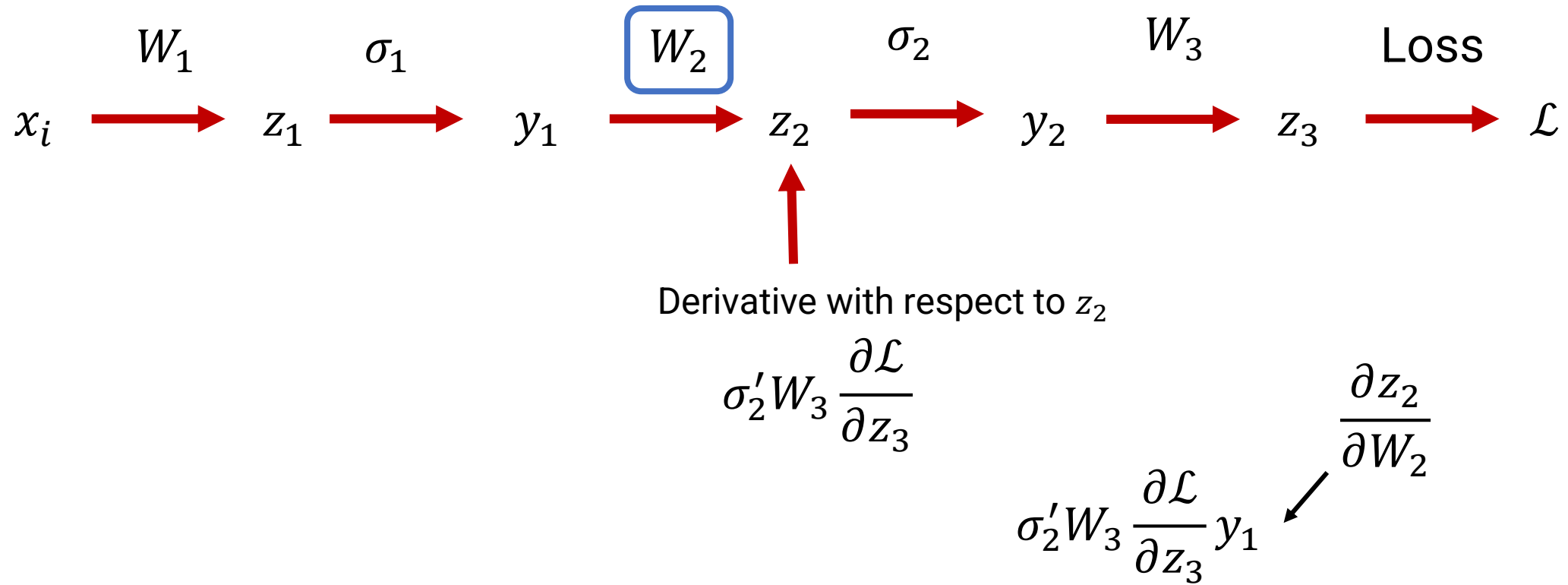
# Derivative With Respect To $W_2$

Derivative of loss  $\mathcal{L}$  with respect to  $W_2$

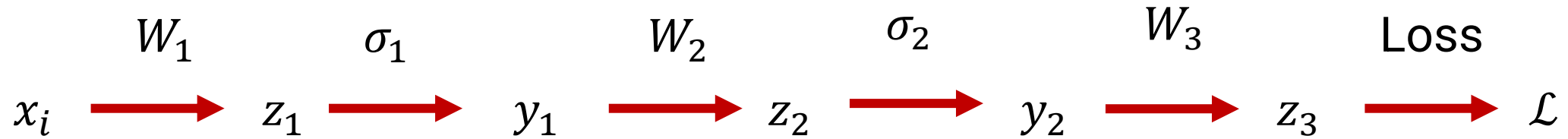


# Derivative With Respect To $W_2$

Derivative of loss  $\mathcal{L}$  with respect to  $W_2$



# Complete Backprop



$$\nabla_{z_3} \mathcal{L}$$

$$\nabla_{W_3} \mathcal{L} = y_2^T \nabla_{z_3} \mathcal{L}$$

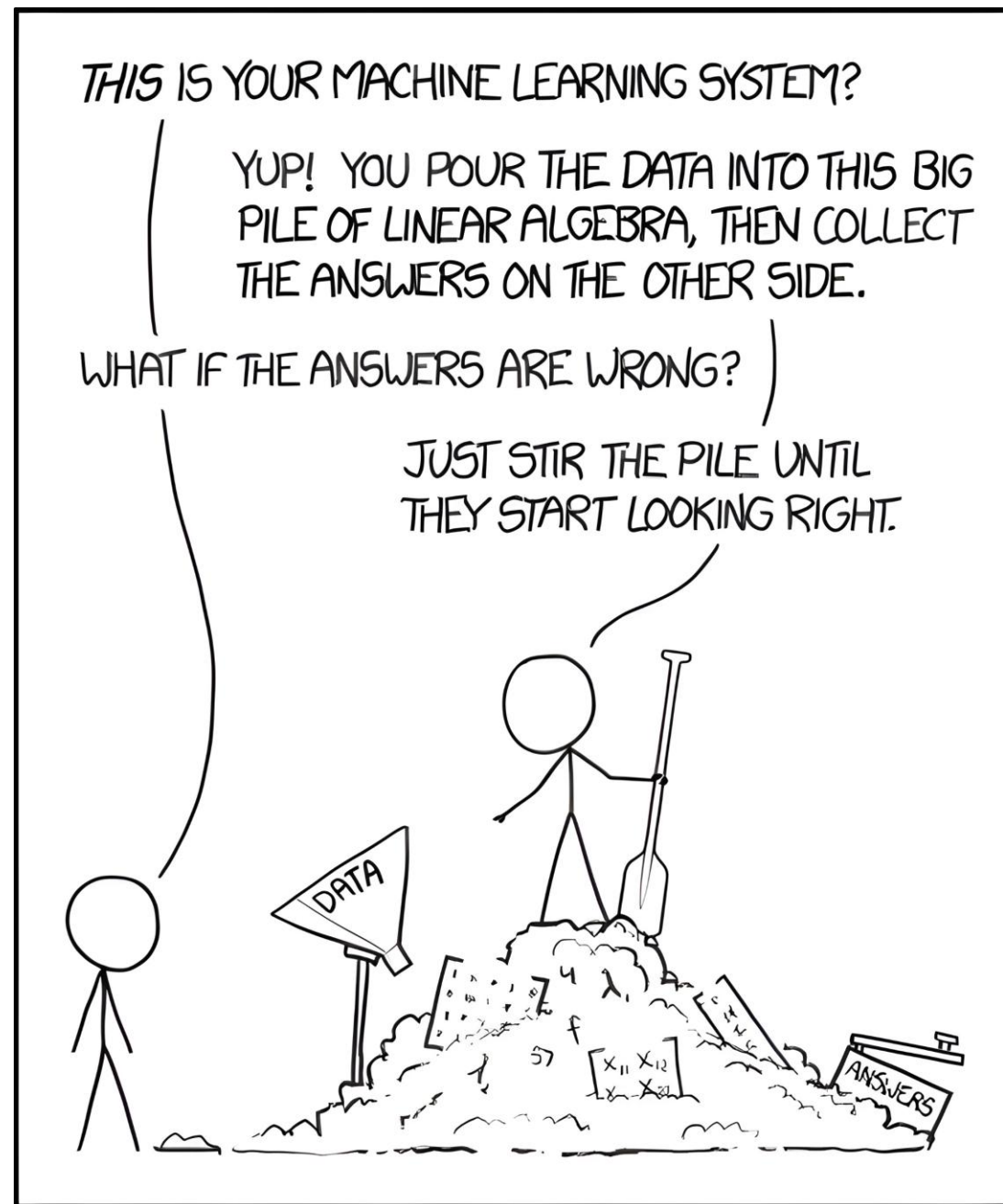
$$\nabla_{z_2} \mathcal{L} = \nabla_{z_3} \mathcal{L} W_3^T \sigma_2'$$

$$\nabla_{W_2} \mathcal{L} = y_1^T \nabla_{z_2} \mathcal{L}$$

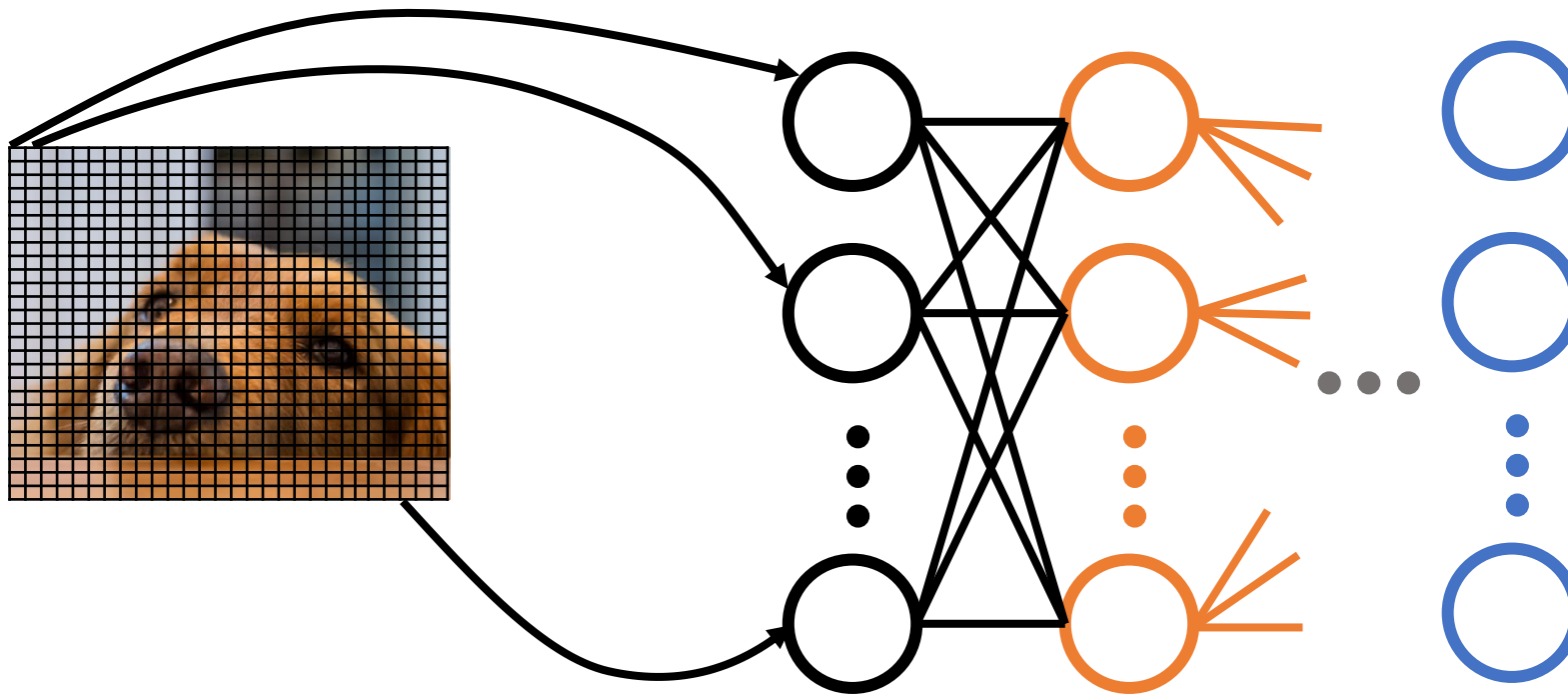
$$\nabla_{z_1} \mathcal{L} = \nabla_{z_3} \mathcal{L} W_3^T \sigma_2' W_2^T \sigma_1'$$

$$\nabla_{W_1} \mathcal{L} = x_i^T \nabla_{z_1} \mathcal{L}$$

# Essentially!



# We Did Something Cool!



Dimension of Output?

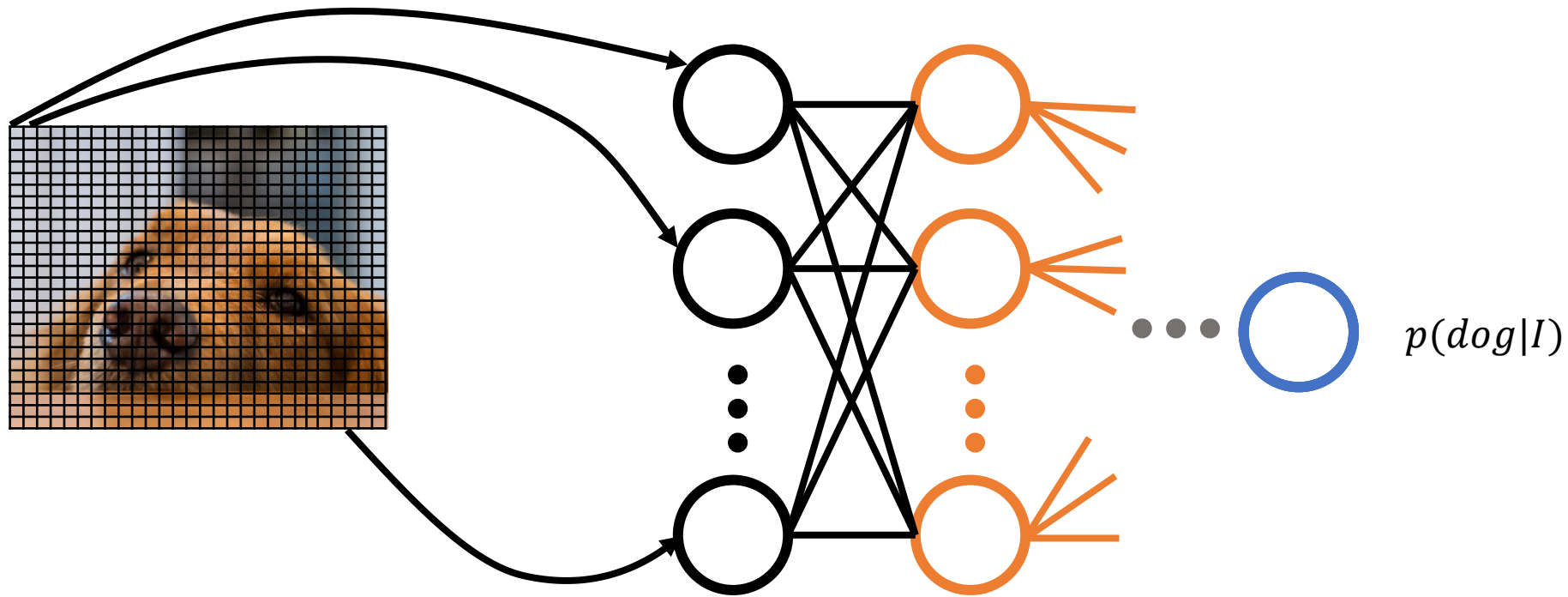
1

Why?

$p(\text{dog}|I)$  as

$p(\text{cat}|I) = 1 - p(\text{dog}|I)$

# We Did Something Cool!



**Fully Connected Network or Multi Layer Perceptron (MLP)**

# Gradient Vs Derivative

$$\frac{\partial \mathcal{L}}{\partial x_1} = W_1 \sigma_1' W_2 \sigma_2' W_3 \frac{\partial \mathcal{L}}{\partial z_3} \quad \text{Derivative}$$

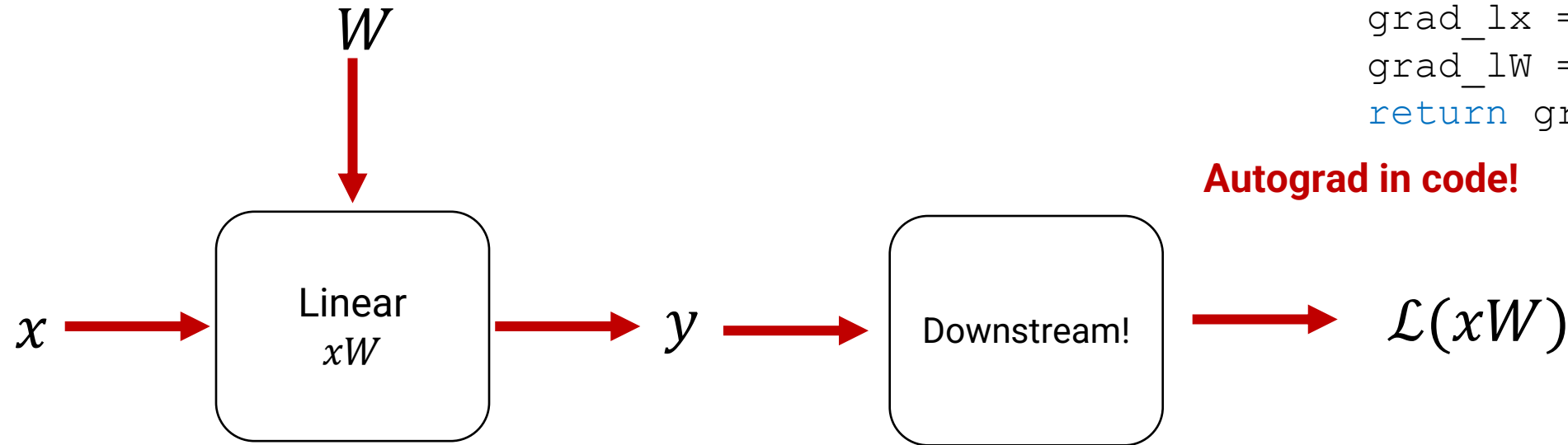
$$\nabla \mathcal{L}_{x_1} = \nabla_{z_3} \mathcal{L} W_3^T \sigma_2' W_2^T \sigma_1' W_1^T \quad \text{Gradient}$$

**Backward Pass!**





# Let's Do A Trick!



```
def dorward(x,W):  
    return x@W
```

```
def backward(x, W, grad_ly):  
    grad_lx = grad_ly@W.T  
    grad_lW = x.T@grad_ly  
    return grad_lx, grad_lW
```

**Autograd in code!**

Gradient with respect to input?

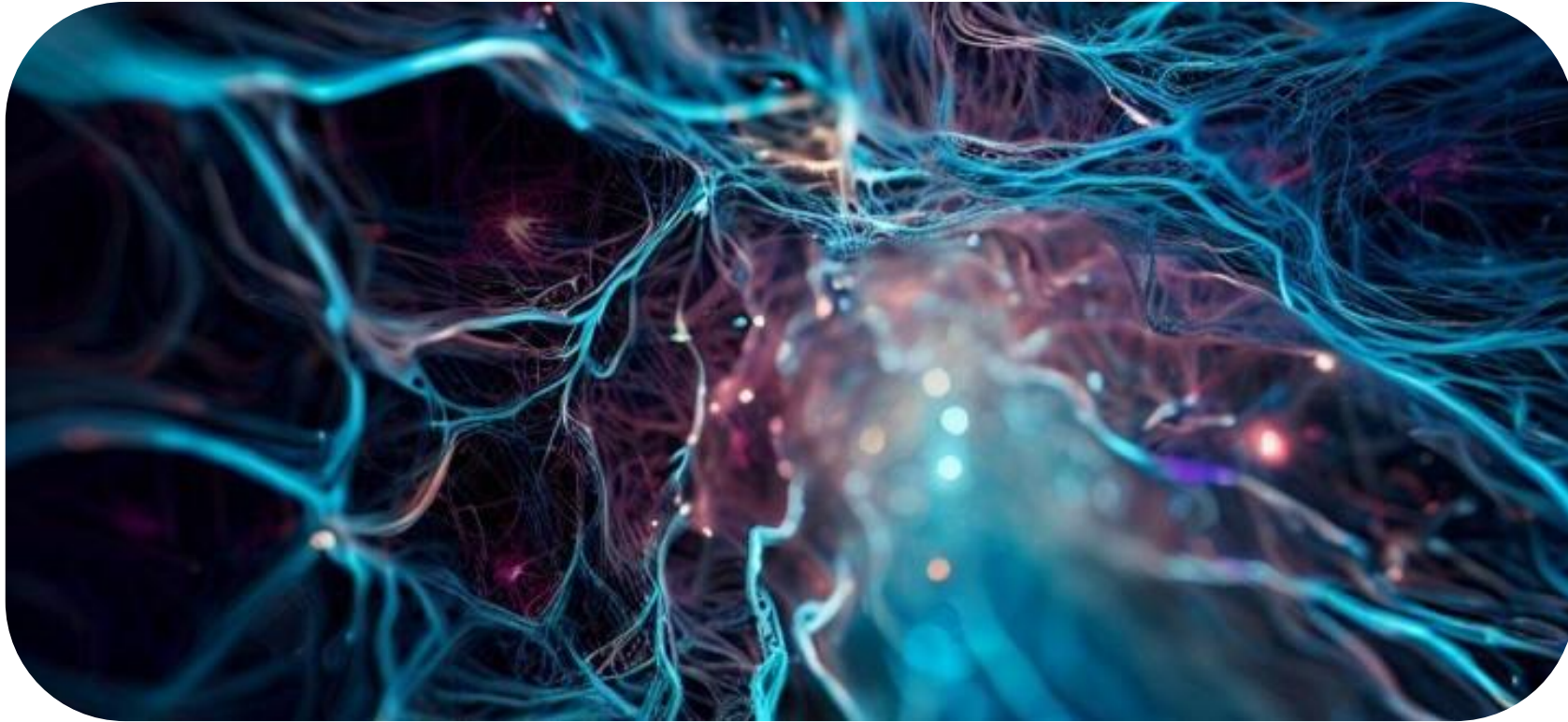
$$\nabla_x \mathcal{L}(xW) = \mathcal{L}'(xW)W^T$$

$$\nabla_W \mathcal{L}(xW) = x^T \mathcal{L}'(xW)$$

Since if  $h(x) = f(Ax) \Rightarrow \nabla h(x) = A^T f'(Ax)$

And if  $h(x) = f(xA) \Rightarrow \nabla h(x) = f'(xA)A^T$

# Next Class!



**NN Tuning, Image Filtering And Convolutional Neural Networks**