# Part 4:
# Deep Neural Networks

● ● ●

Mikhail Romanov

# Gradient Vanishing

# Two Layer Neural Net

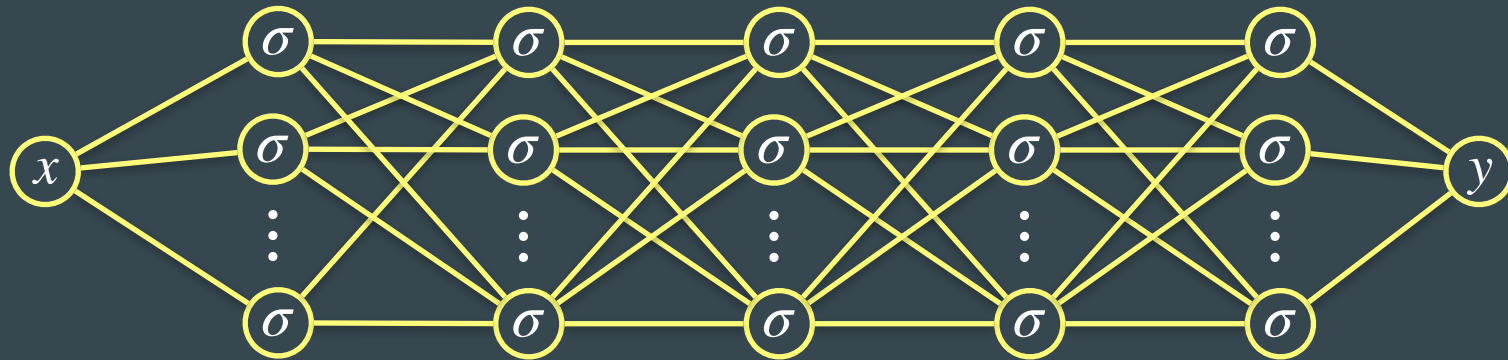$$y = w'_1 \sigma(w_1 x + b_1) + b'_1$$

$$+ w'_2 \sigma(w_2 x + b_2) + b'_2$$
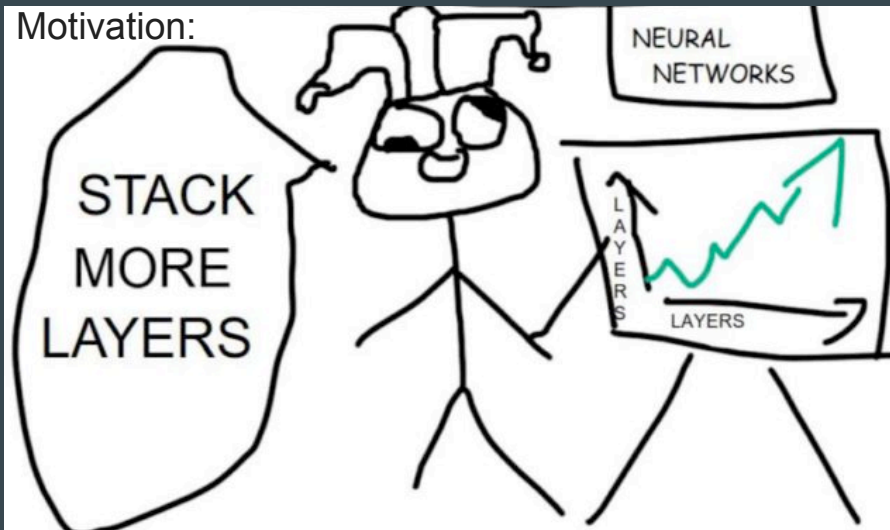
$$+ w'_3 \sigma(w_3 x + b_3) + b'_3$$

$$\hat{y} = b^{out} + \sum_{i=1}^{N} w_i^{out} \sigma(w^i x + b^i)$$



$$\hat{\mathbf{y}} = \mathbf{b}_{out} + \sum_{i=1}^{N} \mathbf{W}_{out} \sigma(\mathbf{W}_{\mathbf{in}} \mathbf{X} + \mathbf{b}_{in})$$

Linear Combination of Sigmoids is
Full System!

# Multi-Layer NNs



$$x$$
$$\text{Linear } 1 \rightarrow N_1$$
$$\text{Sigmoid}$$
$$\text{Linear } N_1 \rightarrow N_2$$
$$\text{Sigmoid}$$
$$\text{Linear } N_2 \rightarrow N_3$$
$$\text{Sigmoid}$$
$$\vdots$$
$$\text{Sigmoid}$$
$$\text{Linear } N_{M-1} \rightarrow N_M$$
$$y$$

Motivation:

# Gradient Vanishing: Sigmoid



$$x$$
$$\text{Linear } 1 \to \text{N}_1$$
$$\text{Sigmoid}$$
$$\text{Linear } \text{N}_1 \to \text{N}_2$$
$$\text{Sigmoid}$$
$$\text{Linear } \text{N}_2 \to \text{N}_3$$
$$\text{Sigmoid}$$
$$\vdots$$
$$\text{Sigmoid}$$
$$\text{Linear } \text{N}_{M-1} \to \text{N}_M$$
$$y$$

$$y = f_1(f_2(f_3(\ldots)))$$

$$\frac{\partial y}{\partial x} = \frac{\partial f_N}{\partial f_{N-1}} \frac{\partial f_{N-1}}{\partial f_{N-2}} \frac{\partial f_{N-2}}{\partial f_{N-3}} \ldots \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial x}$$
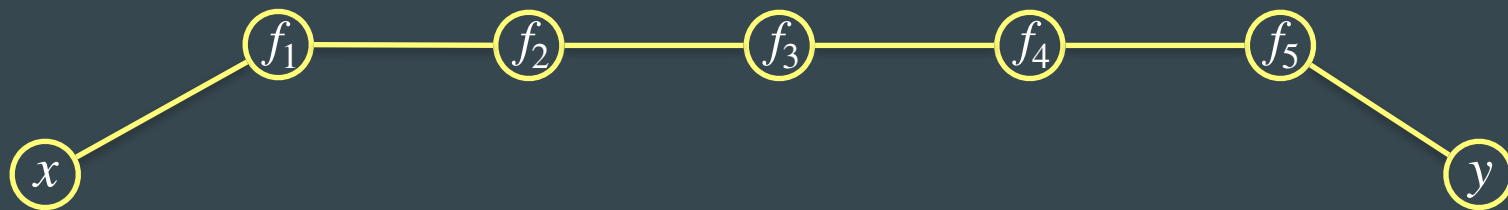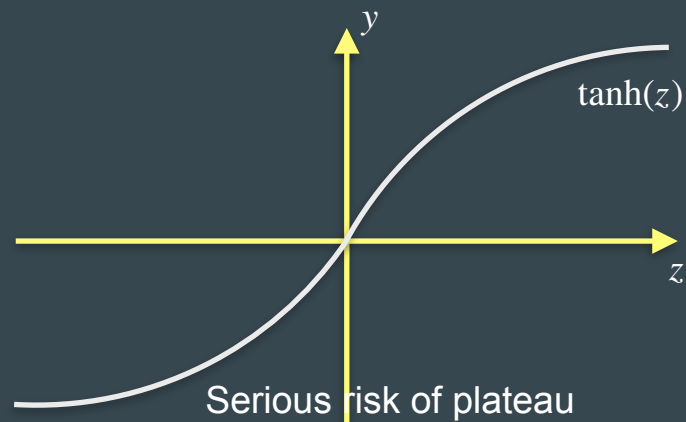
$$\sigma$$

$$\frac{\partial y}{\partial x} = \prod_{i=1}^{N} \frac{\partial \text{Linear}_i}{\partial z_i} \prod_{j=1}^{M} \frac{\partial \sigma_j}{\partial z_j} = V \prod_{j=1}^{M} \sigma(z_j)(1 - \sigma(z_j)) \leq V \frac{1}{4^M}$$

Serious risk of plateau

# Tanh

$f_1$ — $f_2$ — $f_3$ — $f_4$ — $f_5$

$x$ ⟋ ⟍ $y$

$$y = f_1(f_2(f_3(\dots)))$$

$x$
Linear $1 \to N_1$
Tanh
Linear $N_1 \to N_2$
Tanh
Linear $N_2 \to N_3$
Tanh
⋮
Tanh
Linear $N_{M-1} \to N_M$
$y$

$$\frac{\partial y}{\partial x} = \frac{\partial f_N}{\partial f_{N-1}} \frac{\partial f_{N-1}}{\partial f_{N-2}} \frac{\partial f_{N-2}}{\partial f_{N-3}} \dots \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial x}$$

$$\tilde{\sigma} = \tanh$$

$$\frac{\partial y}{\partial x} = \prod_{i=1}^{N} \frac{\partial \text{Linear}_i}{\partial z_i} \prod_{j=1}^{M} \frac{\partial \tilde{\sigma}_j}{\partial z_j} = V \prod_{j=1}^{M} (1 + \tilde{\sigma}(z_j))(1 - \tilde{\sigma}(z_j)) \leq V 1^M$$

$y$

$\tanh(z)$

$z$

Serious risk of plateau

# Leaky ReLU



$x$
Linear $1 \rightarrow N_1$
L $-$ ReLU
Linear $N_1 \rightarrow N_2$
L $-$ ReLU
Linear $N_2 \rightarrow N_3$
L $-$ ReLU
$\vdots$
L $-$ ReLU
Linear $N_{M-1} \rightarrow N_M$
$y$
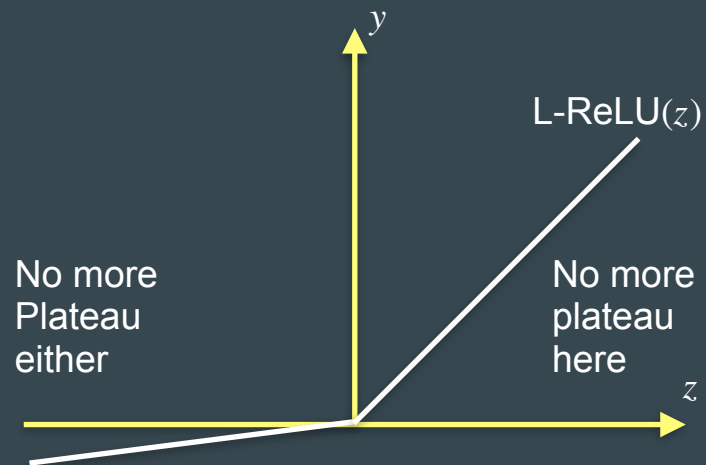
$y = f_1(f_2(f_3(\dots)))$

$$\frac{\partial y}{\partial x} = \frac{\partial f_N}{\partial f_{N-1}} \frac{\partial f_{N-1}}{\partial f_{N-2}} \frac{\partial f_{N-2}}{\partial f_{N-3}} \dots \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial x}$$

L-ReLU($z$)

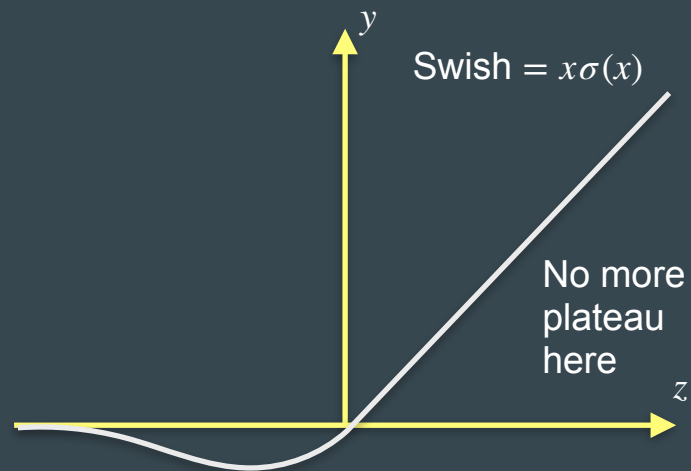No more
Plateau
either

No more
plateau
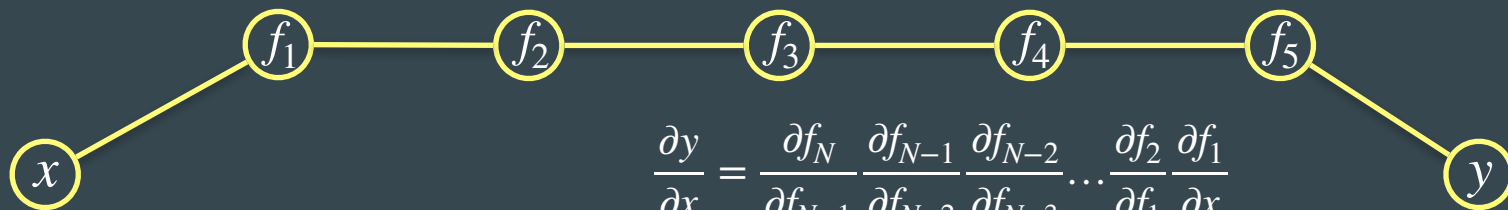here

# Leaky ReLU



$$y = f_1(f_2(f_3(\dots)))$$

$$\frac{\partial y}{\partial x} = \frac{\partial f_N}{\partial f_{N-1}} \frac{\partial f_{N-1}}{\partial f_{N-2}} \frac{\partial f_{N-2}}{\partial f_{N-3}} \dots \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial x}$$

$x$
Linear $1 \to N_1$
L − ReLU
Linear $N_1 \to N_2$
L − ReLU
Linear $N_2 \to N_3$
L − ReLU
$\vdots$
L − ReLU
Linear $N_{M-1} \to N_M$
$y$

Swish $= x\sigma(x)$

No more plateau here

# Linear Layer



$$x$$
$$\text{Linear } 1 \rightarrow N_1$$
$$L - ReLU$$
$$\text{Linear } N_1 \rightarrow N_2$$
$$L - ReLU$$
$$\text{Linear } N_2 \rightarrow N_3$$
$$L - ReLU$$
$$\vdots$$
$$L - ReLU$$
$$\text{Linear } N_{M-1} \rightarrow N_M$$
$$y$$

$$y = f_1(f_2(f_3(\dots)))$$

$$\frac{\partial y}{\partial x} = \frac{\partial f_N}{\partial f_{N-1}} \frac{\partial f_{N-1}}{\partial f_{N-2}} \frac{\partial f_{N-2}}{\partial f_{N-3}} \dots \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial x}$$

$$z = ax + b$$

$$\partial_x z = a$$

# What about linear operations?



$$\frac{\partial y}{\partial x} = \frac{\partial f_N}{\partial f_{N-1}} \frac{\partial f_{N-1}}{\partial f_{N-2}} \frac{\partial f_{N-2}}{\partial f_{N-3}} \ldots \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial x}$$

$x$
Linear $1 \to N_1$
Sigmoid
Linear $N_1 \to N_2$
Sigmoid
Linear $N_2 \to N_3$
Sigmoid
$\vdots$
Sigmoid
Linear $N_{M-1} \to N_M$
$y$

$$y = f_1(f_2(f_3(\ldots)))$$

$$z = wx + b \qquad\qquad \mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

$$\frac{\partial L}{\partial x} = w \frac{\partial L}{\partial z} \qquad\qquad \frac{\partial L}{\partial \mathbf{x}} = \mathbf{W}^T \frac{\partial L}{\partial \mathbf{z}}$$

One cannot avoid gradient Vanishing!
AT ALL!

$$\left| \frac{\partial L}{\partial x} \right| = w \left| \frac{\partial L}{\partial z} \right| \qquad\qquad \sigma_{min} \left| \frac{\partial L}{\partial z} \right| \leq \left| \frac{\partial L}{\partial x} \right| \leq \sigma_{max} \left| \frac{\partial L}{\partial z} \right|$$

# Can we avoid gradient vanishing?



$$y = f_1(f_2(f_3(\ldots)))$$

$$\frac{\partial y}{\partial x} = \frac{\partial f_N}{\partial f_{N-1}} \frac{\partial f_{N-1}}{\partial f_{N-2}} \frac{\partial f_{N-2}}{\partial f_{N-3}} \cdots \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial x}$$

$x$
Linear $1 \rightarrow N_1$
Sigmoid
Linear $N_1 \rightarrow N_2$
Sigmoid
Linear $N_2 \rightarrow N_3$
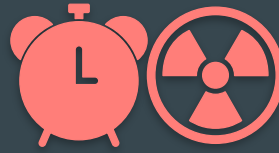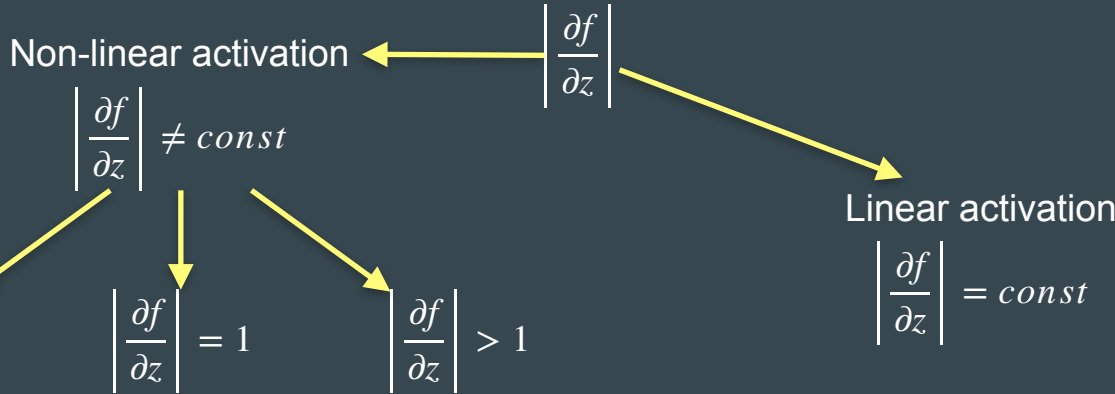Sigmoid
$\vdots$
Sigmoid
Linear $N_{M-1} \rightarrow N_M$
$y$

Non-linear activation $\longleftarrow$ $\left| \dfrac{\partial f}{\partial z} \right|$

$\left| \dfrac{\partial f}{\partial z} \right| \neq const$

Linear activation

$\left| \dfrac{\partial f}{\partial z} \right| = const$

$\left| \dfrac{\partial f}{\partial z} \right| < 1$     $\left| \dfrac{\partial f}{\partial z} \right| = 1$     $\left| \dfrac{\partial f}{\partial z} \right| > 1$

One cannot avoid gradient Vanishing!

# Residual Connections

# Residual Connection



$$y = f_1(f_2(f_3(\ldots)))$$

$$\partial_x y = \frac{\partial f_N}{\partial z_{N-1}} \frac{\partial f_{N-1}}{\partial z_{N-2}} \frac{\partial f_{N-2}}{\partial z_{N-3}} \cdots \frac{\partial f_2}{\partial z_1} \frac{\partial f_1}{\partial x}$$

$$z = f(x) + x$$

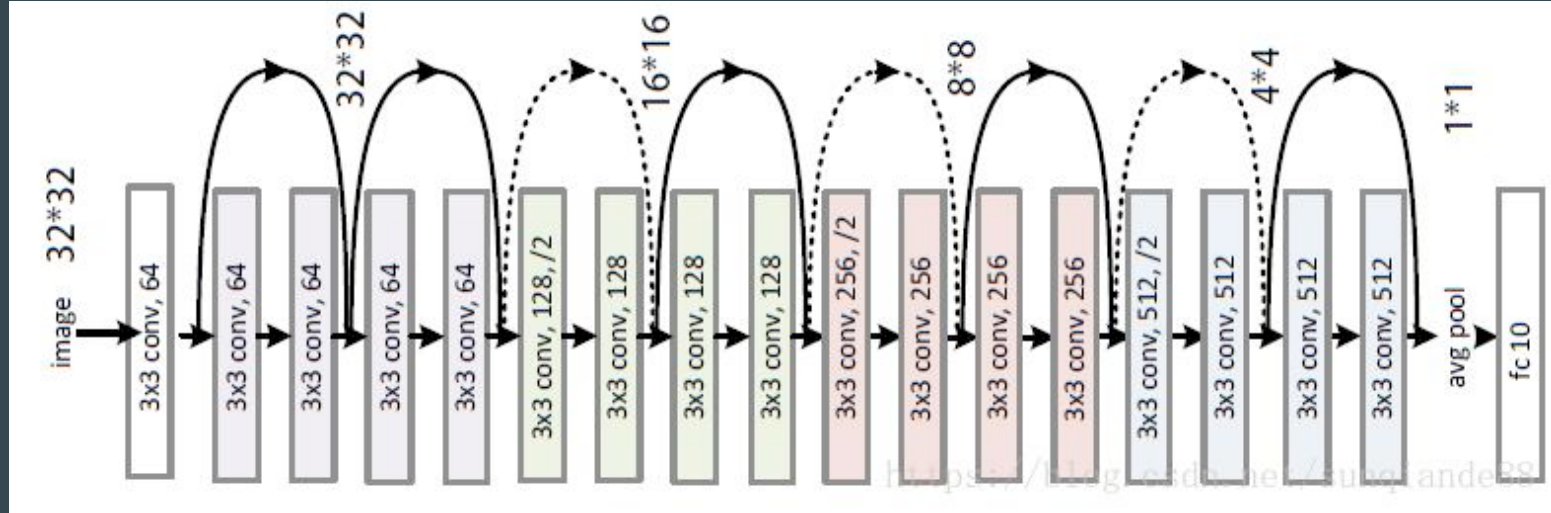$$\partial_x z = \partial_x f(x) + 1$$

$$y = f_1(f_2(f_3(\ldots)))$$

$$\frac{\partial y}{\partial x} = \left(1 + \frac{\partial f_N}{\partial z_{N-1}}\right) \cdots \left(1 + \frac{\partial f_2}{\partial z_1}\right) \left(1 + \frac{\partial f_1}{\partial x}\right)$$

# ResNet Allows: Extremely Deep Networks

# ResNet Allows: "Boosted" Neural Networks



Network-in-Network

Linear N -> M

$\sigma$

Linear M -> N

+

# Normal signals

# What if the signals are normal



Gradient amplitudes:

Bad case (

Very bad case (((

OK

# How to preserve normality

$x:$

$x_1$

$\mathbb{E}x = 0$

$\mathbb{V}x = 1$

$x_2$

$x_3$

$w_1$

$w_2$

$w_3$

$\Sigma$

$y$

$y = \displaystyle\sum_{i=1}^{N} w_i x_i$

$y \propto ?$

$\mathbb{E}y = \mathbb{E}\left( \displaystyle\sum_{i=1}^{N} w_i x_i \right) = \displaystyle\sum_{i=1}^{N} w_i \mathbb{E}x_i = 0$

$\mathbb{V}y = \mathbb{V}\left( \displaystyle\sum_{i=1}^{N} w_i x_i \right) = \displaystyle\sum_{i=1}^{N} \mathbb{V}\left( w_i x_i \right) = \displaystyle\sum_{i=1}^{N} w_i^2 \mathbb{V}x = \displaystyle\sum_{i=1}^{N} w_i^2$

Zero mean is preserved

How to keep unit variance?

$\displaystyle\sum_{i=1}^{N} w_i^2 \approx 1$

$w_i \propto \mathcal{N}\left( 0, \dfrac{1}{\sqrt{N}} \right)$     He initialisation

$w_i \propto \mathcal{U}\left( -\dfrac{C}{\sqrt{N}}, \dfrac{C}{\sqrt{N}} \right)$ Xavier initialisation

# How to enforce normality on an input

$$\mathbf{x}^* = \frac{\mathbf{x} - \mu}{\sigma}$$

$$\mu = \frac{1}{N} \sum_{s=1}^{S} \mathbf{x}_s$$

$$\sigma = \sqrt{\frac{1}{N-1} \sum_{s=1}^{S} (\mathbf{x}_s - \mu)^2}$$

1) We set mean to zero

2) We set standard deviation to one

Now the inputs are perfectly OK!

5-Sigma rule:
- Now the input signals are bound
- To the interval [-5, 5]
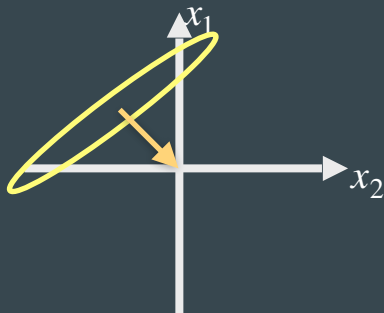- And only 1 out of 1 million
- Leaves this interval

# What the normality is

$$\mathbf{x}^* = \Sigma^{-1}(\mathbf{x} - \mu)$$

$$\mu = \frac{1}{N}\sum_{s=1}^{S}\mathbf{x}_s$$

$$\Sigma^T\Sigma = \frac{1}{N-1}\sum_{s=1}^{S}(\mathbf{x}_s - \mu)(\mathbf{x}_s - \mu)^T$$



Centering
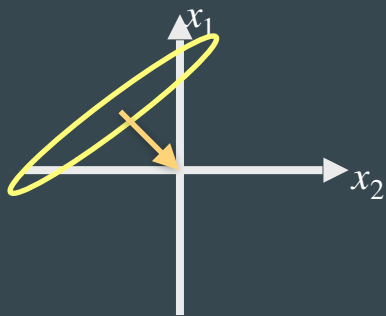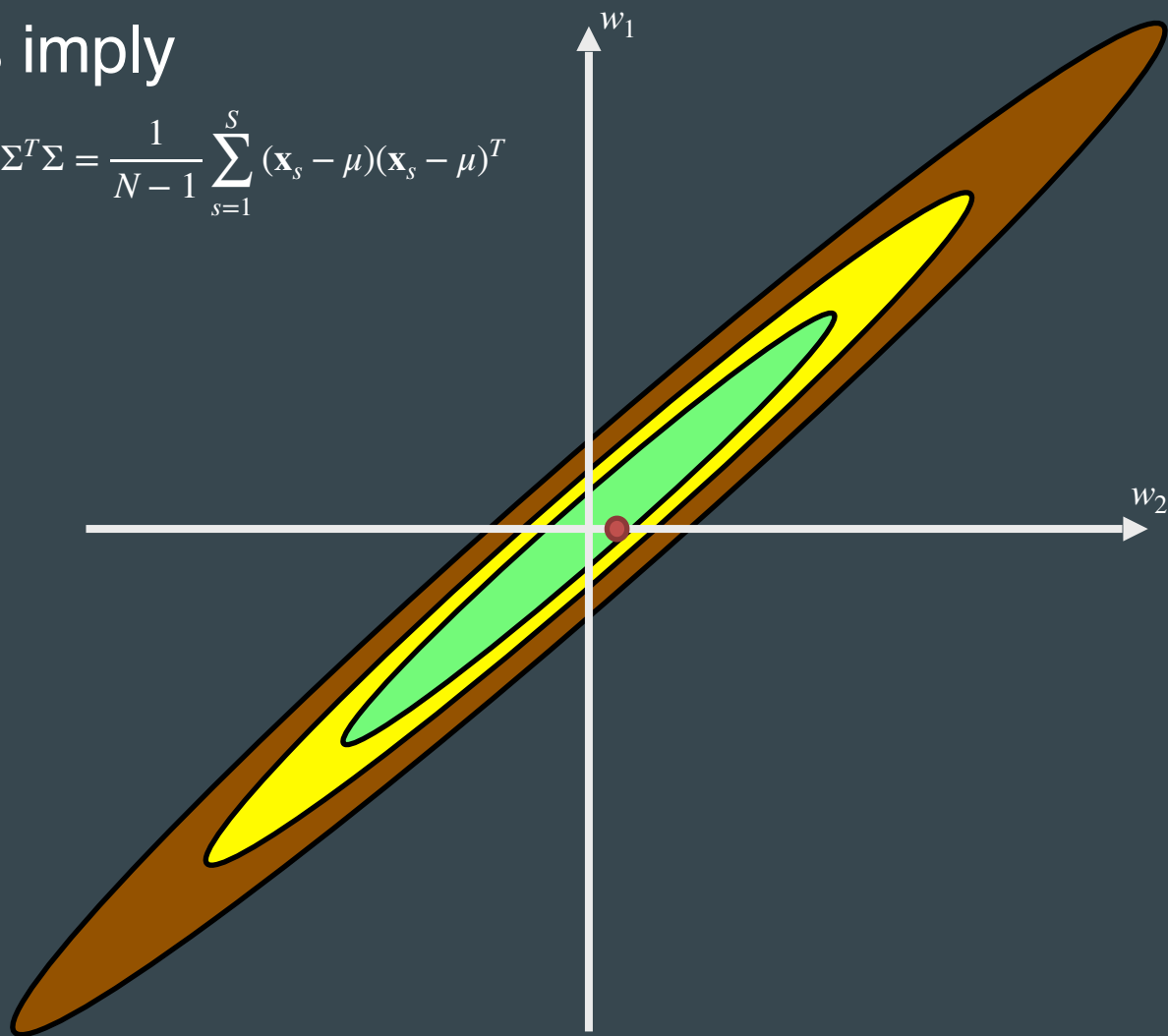
Decorrelation

Standardisation

Result

# What each of the steps imply

$$\mu = \frac{1}{N} \sum_{s=1}^{S} \mathbf{x}_s$$

$$\Sigma^T \Sigma = \frac{1}{N-1} \sum_{s=1}^{S} (\mathbf{x}_s - \mu)(\mathbf{x}_s - \mu)^T$$

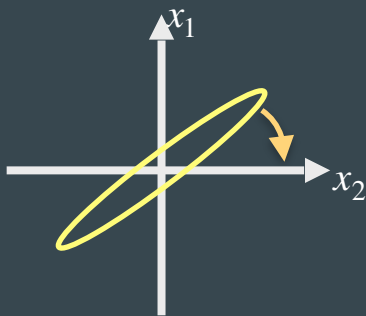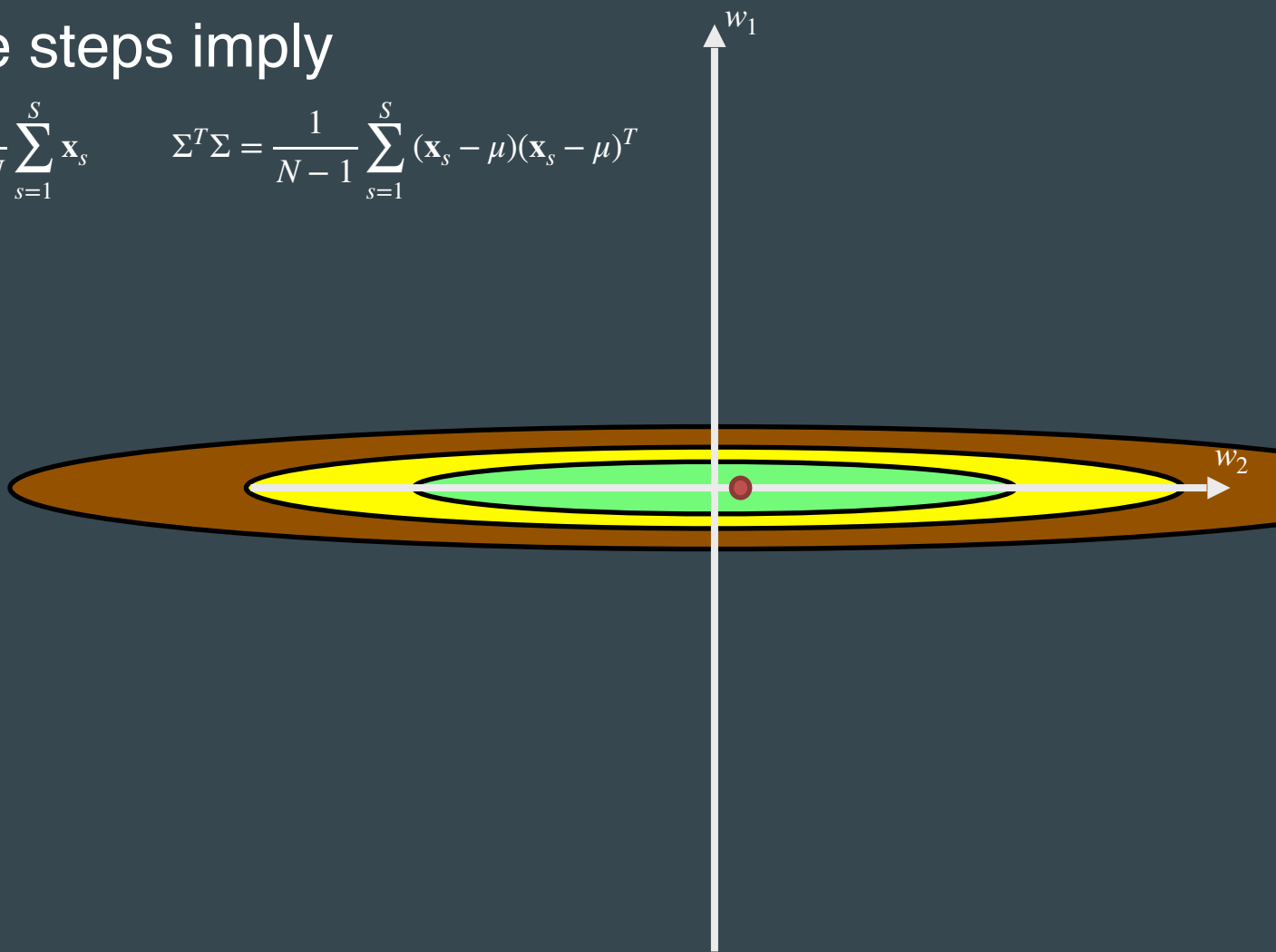$$\mathbf{x}^* = \Sigma^{-1}(\mathbf{x} - \mu)$$

Centering

# What each of the steps imply

$$\mathbf{x}^* = \Sigma^{-1}(\mathbf{x} - \mu) \qquad \mu = \frac{1}{N}\sum_{s=1}^{S}\mathbf{x}_s \qquad \Sigma^T\Sigma = \frac{1}{N-1}\sum_{s=1}^{S}(\mathbf{x}_s - \mu)(\mathbf{x}_s - \mu)^T$$
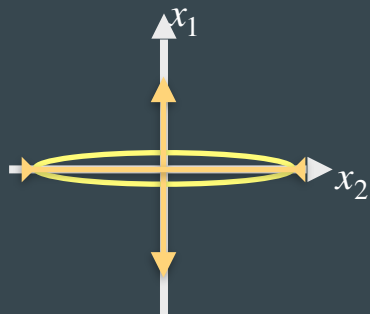
Decorrelation

# What each of the steps imply

$$\mathbf{x}^* = \Sigma^{-1}(\mathbf{x} - \mu) \qquad \mu = \frac{1}{N}\sum_{s=1}^{S}\mathbf{x}_s \qquad \Sigma^T\Sigma = \frac{1}{N-1}\sum_{s=1}^{S}(\mathbf{x}_s - \mu)(\mathbf{x}_s - \mu)^T$$
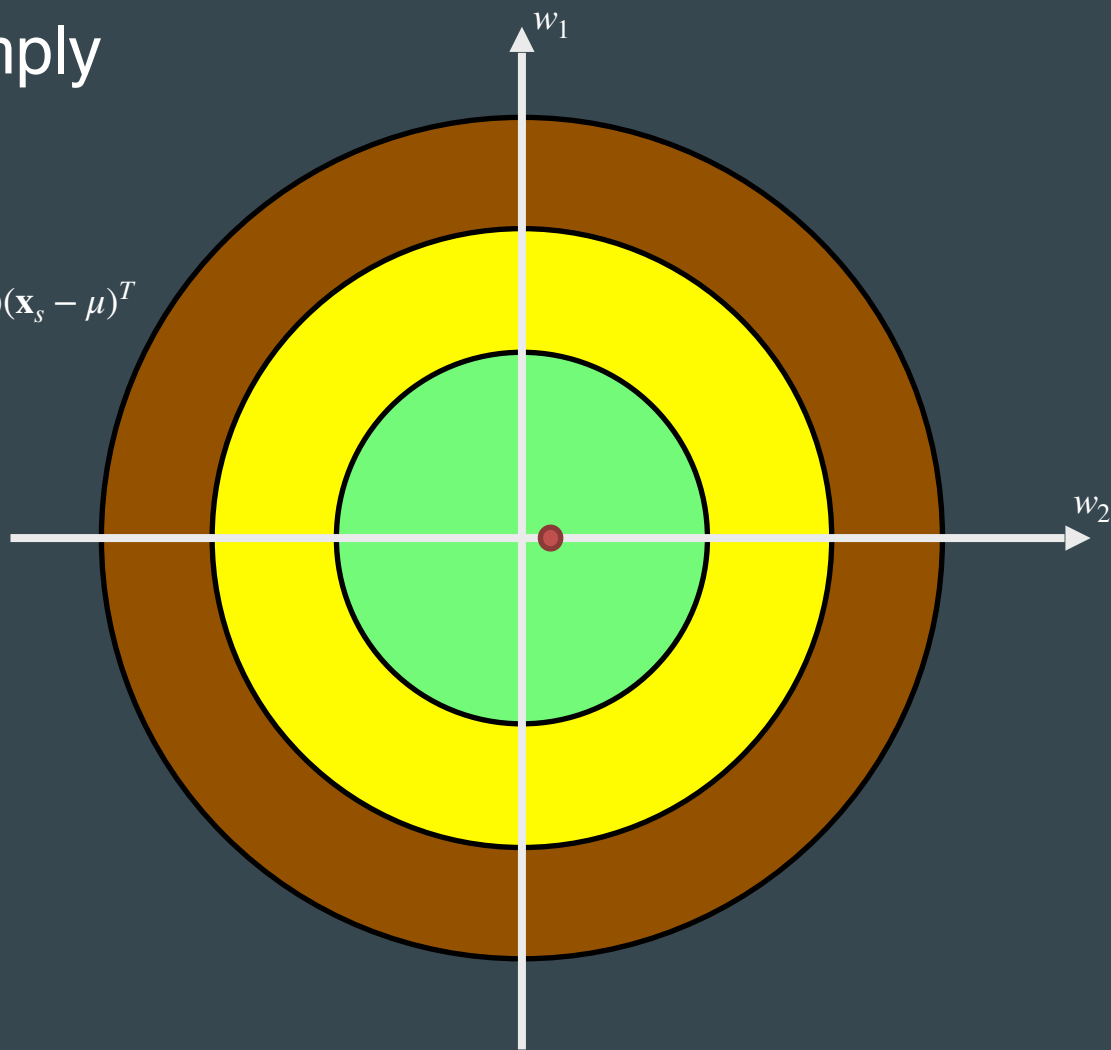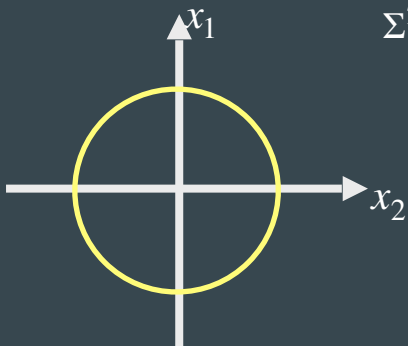
Standardisation

# What each of the steps imply
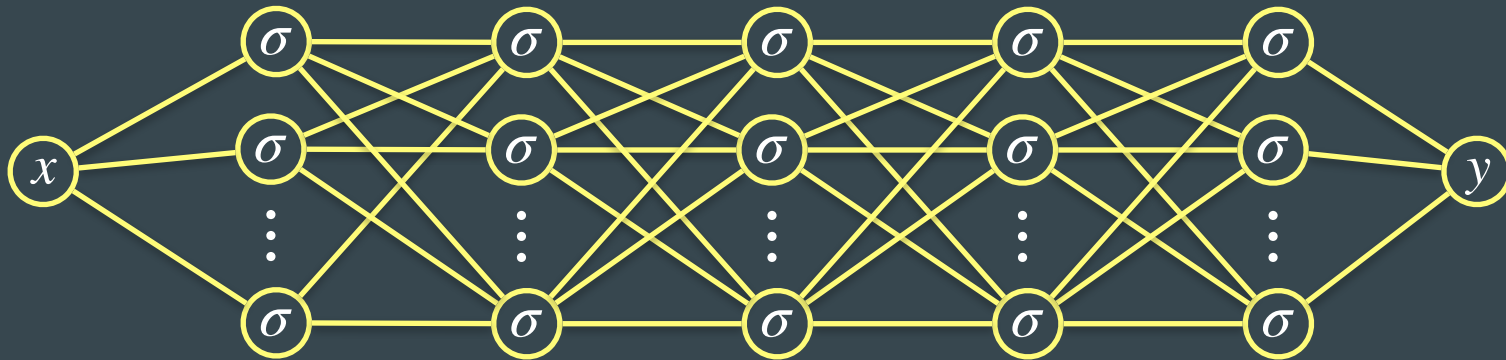
$$\mathbf{x}* = \Sigma^{-1}(\mathbf{x} - \mu)$$

$$\mu = \frac{1}{N}\sum_{s=1}^{S}\mathbf{x}_s$$

Result

$$\Sigma^T\Sigma = \frac{1}{N-1}\sum_{s=1}^{S}(\mathbf{x}_s - \mu)(\mathbf{x}_s - \mu)^T$$

# How to enforce normality on signals



$$z* = \frac{\mathbf{z} - \mu}{\sigma}\mathbf{a} + \mathbf{b}$$

| | | |
|---|---|---|
| $\mu, \sigma$ | Statistical parameters | Computed on one batch |
| $a, b$ | Trainable parameters | Optimised |

Validation:

$$\hat{\mu} = EMA(\mu)$$

$$\hat{\sigma}^2 = EMA(\sigma^2)$$

$z$

# Other reasons why that is useful

# Regularisation

# What the regularisation is

$$Ax = b$$

$$x = A^{-1}b$$

$$\|Ax - b\|_2^2 = 0$$

$$x = \underset{x}{\mathrm{argmax}}\|Ax - b\|_2^2$$

One solution if

$$\dim(x) \leq \dim(b)$$

No issues here

Infinitely many solutions if

$$\dim(x) > \dim(b)$$

We are interested only in the simpliest solution

How to measure the simplicity of a solution?

$$\|x\| \quad \text{— complexity}$$
$$\text{measure}$$

$$L = \|Ax - b\|_2^2$$

Does not take into account solution complexity

Many solutions

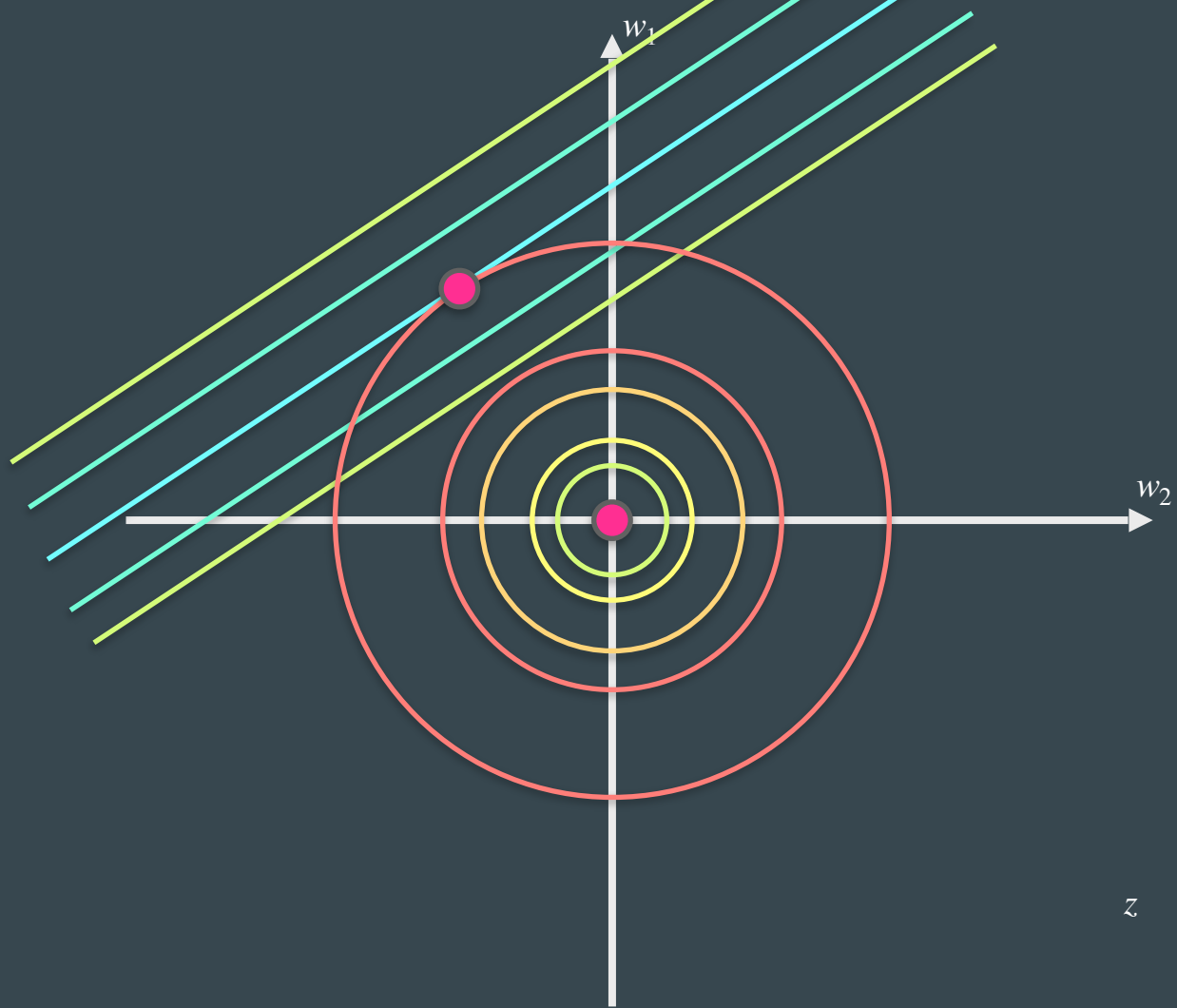$$L^* = \|Ax - b\|_2^2 + \lambda\|x\|_2^2$$

Takes into account solution complexity

But displaces the solution

Only one solution
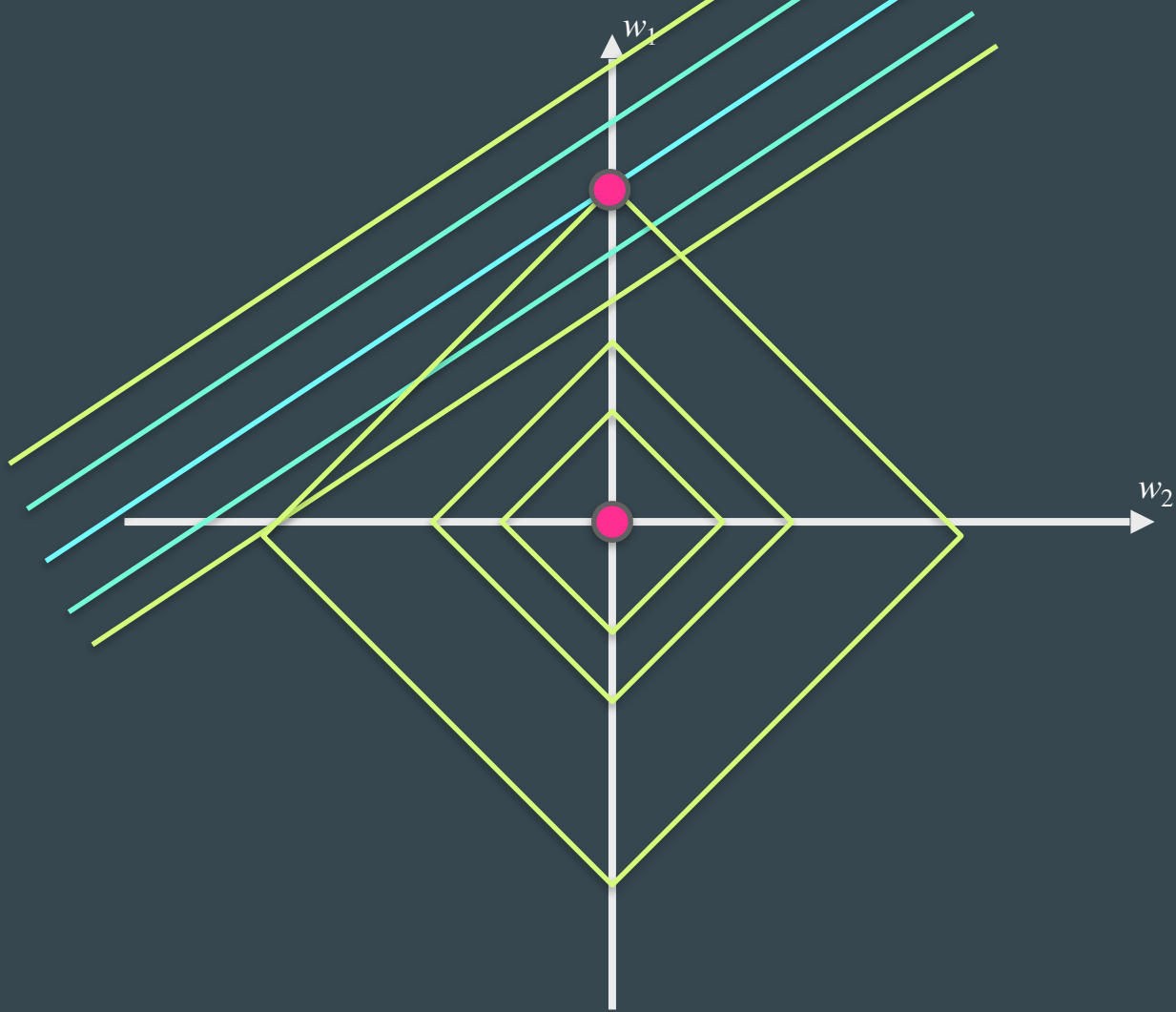
$z$

# L2 regularisation

$$L* = \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda\|\mathbf{x}\|_2^2$$

# L1 regularisation

$L* = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda\|\mathbf{x}\|_M$

$\|\mathbf{x}\|_M = x_1 + x_2 + \ldots + x_N$

# Summary

- Gradient Vanishing problem and its source
- Choice: Gradient Vanishing or Gradient Explosion
- Methods of mitigating Gradient Vanishing
- Batch Normalisation: enforcing the Normality on Neural Network's signals
- Other positive sides of BatchNorms
- Regularisation and why is that good