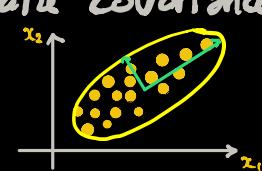
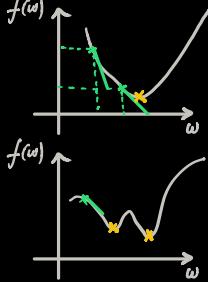


# M&L<sub>4</sub>

1. Recap on date covariance  $\Sigma$
- 
- data shift  
linear & non linear Bayes classifier

- 2 - Fundamentals of numerical optimization  $[\min_{\omega} l(\omega)]$

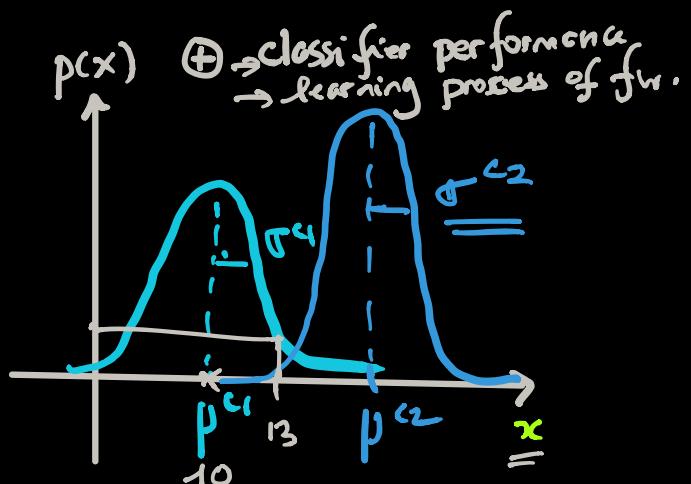
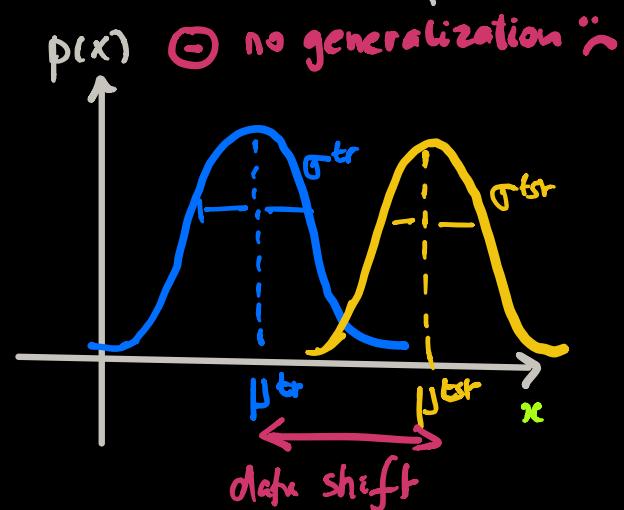


- └ Taxonomy of supervised learning
- └ Taylor approximation of a function
- └ Stationary points
- └ Convex and non-convex functions
- └ gradient descent

# MööL<sub>3</sub>?

$X = N \begin{bmatrix} \text{tr} \\ \vdots \\ \text{tst} \end{bmatrix}$

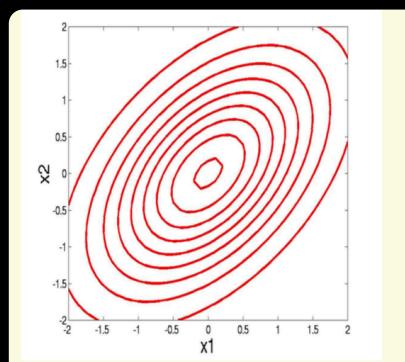
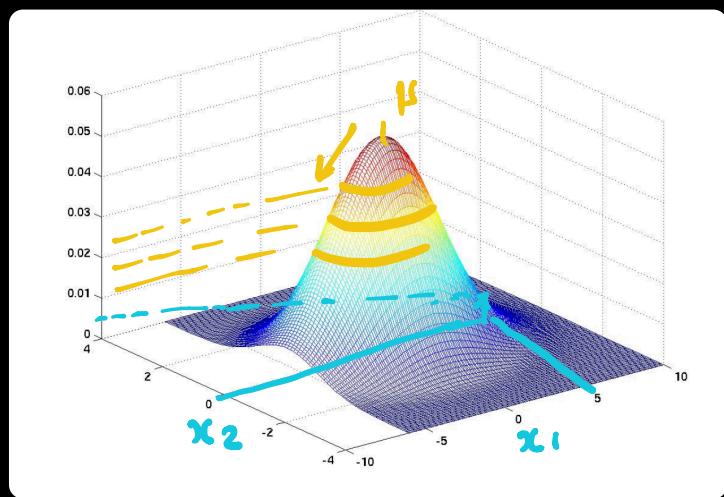
$f_w \rightarrow y = \begin{bmatrix} +1 \\ -1 \\ +1 \\ -1 \\ +1 \\ +1 \end{bmatrix}$



$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right]$$

$\underline{p(x=13)}$  ← the probability of feature  $x$  taking value 13?

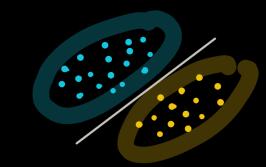
$x \in \mathbb{R}^d$ , multivariate distribution



$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

$$\underline{p(x)} \in \mathbb{R}^d = p(x_1, x_2, \dots, x_d) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right]$$

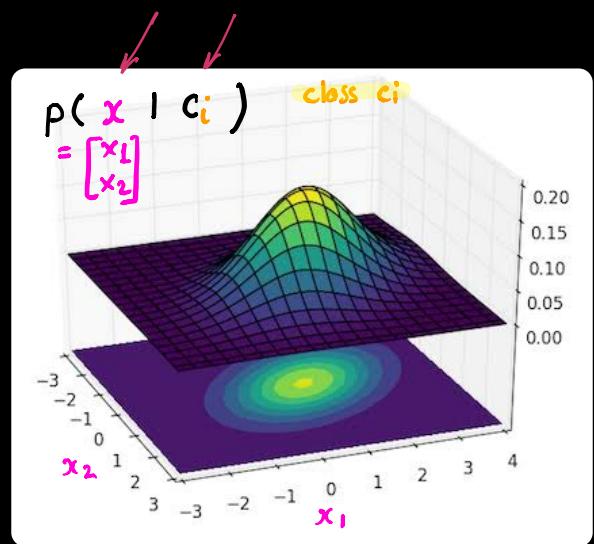
# Bayes discriminant classifier



- new testing sample  $\mathbf{x}$
- what is the probability of  $\mathbf{x}$  belonging to class  $c_i$ ?

$$g_i(\mathbf{x}) = \underbrace{P(c_i | \mathbf{x})}_{\substack{\text{conditional} \\ \text{probability}}} = \frac{P(c_i) P(\mathbf{x} | c_i)}{P(\mathbf{x})} = \frac{\text{likelihood of observing } \mathbf{x} \text{ given } c_i \text{ class}}{\text{marginal likelihood or probability of observing } \mathbf{x}}$$

$\Rightarrow$  we want to predict "class  $c_i$ "

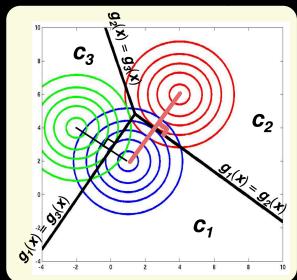


$$\Sigma_i = \sigma^2 I$$

(different means  $\mu_i$  but equal variances)

$$g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0}$$

(linear in  $\mathbf{x}$ )



lines connecting the classes means are  $\perp$  to the decision boundaries

Special case  
 $\ln p(c_i) = \ln p(\mathbf{x}) + \ln p(c_i)$

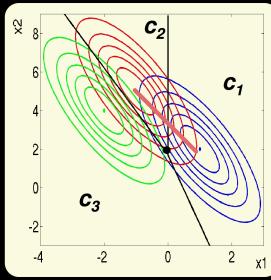
$$g_i(\mathbf{x}) = -\frac{1}{2} \|\mathbf{x} - \mu_i\|_2^2$$

$$\Sigma_i = \Sigma_j = \Sigma$$

(arbitrary covariance but constant across classes)

$$g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0}$$

(linear in  $\mathbf{x}$ )



lines connecting the classes means are  $\not\perp$  to the decision boundaries

Special case  
 $\ln p(c_i) = \ln p(\mathbf{x}) + \ln p(c_i)$

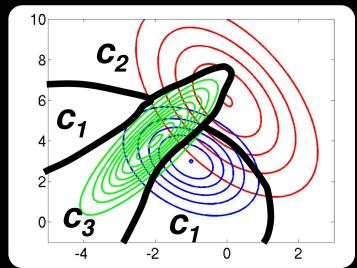
$$g_i(\mathbf{x}) = -\frac{1}{2} \|\mathbf{x} - \mu_i\|_2^2$$

$$(\Sigma_1, \dots, \Sigma_k)$$

( $k$  covariance matrices, general case)

$$g_i(\mathbf{x}) = \underline{\mathbf{x}^t \Sigma^{-1} \mathbf{x}} + \underline{w_i^t \mathbf{x}} + w_{i0}$$

(quadratic in  $\mathbf{x}$ )



decision boundaries between classes are non-linear.

## Mathematical notation

## Definition

$\mathcal{D}$	dataset
$n$	number of samples in a dataset $\mathcal{D}$
$d$	number of features
$\mathbf{x} \in \mathbb{R}^{d \times 1}$	feature vector or data point (sample)
$\mathbf{x}^{\text{sample}}_i \in \mathbb{R}^{d \times 1}$	$i^{\text{th}}$ sample in the population
$\mathbf{x}^{\text{feature}}_j \in \mathbb{R}^{d \times 1}$	$j^{\text{th}}$ feature of $i^{\text{th}}$ sample in the population
$\Sigma \in \mathbb{R}^{d \times d}$	covariance matrix of data population $\{\mathbf{x}^i\}_{i=1}^n$
$ \mathbf{A}  \in \mathbb{R}$	determinant of matrix $A$
$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{\ x - \mu\ _2^2}{\sigma^2}\right)$	probability density function of a variable $x \in \mathbb{R}$
$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}  \Sigma ^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right]$	probability density function of a multidimensional variable $\mathbf{x} \in \mathbb{R}^{N \times 1}$
$\boldsymbol{\mu} \in \mathbb{R}^{d \times 1}$	sample mean $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^i$
$\ \mathbf{x} - \boldsymbol{\mu}\ _{\Sigma^{-1}} = (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \in \mathbb{R}$	Mahalanobis distance between $\mathbf{x}$ and $\boldsymbol{\mu}$
$\mathbf{I}_{d \times d} \in \mathbb{R}^{d \times d}$	identify matrix of size $d \times d$
$\ \mathbf{x} - \boldsymbol{\mu}\ _{\mathbf{I}_{d \times d}} = (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{x} - \boldsymbol{\mu}) \in \mathbb{R}$	Euclidean distance between $\mathbf{x}$ and $\boldsymbol{\mu}$
$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$	also noted as $L_2$ norm $\ \cdot\ _2$
$(i.e., \text{different means } \boldsymbol{\mu}_i \text{ and equal data variances across classes})$	discriminant Bayes function for class $i$ when $\Sigma_i = \sigma^2 \mathbf{I}$
$(i.e., \text{lines connecting means of different classes are perpendicular to decision boundaries})$	if $\ln(p(c_i)) = \ln(p(c_j))$ , $g_i(\mathbf{x}) = -\ \mathbf{x} - \boldsymbol{\mu}_i\ _2^2$
$(i.e., \text{constant data covariance } \Sigma \text{ across classes})$	discriminant Bayes function for class $i$ when $\Sigma_i = \Sigma_j = \Sigma$
$(i.e., \text{lines connecting means of different classes are not perpendicular to decision boundaries})$	if $\ln(p(c_i)) = \ln(p(c_j))$ , $g_i(\mathbf{x}) = -\frac{1}{2} \ \mathbf{x} - \boldsymbol{\mu}_i\ _{\Sigma^{-1}}^2$
$g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W} \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$	quadratic discriminant function (decision boundaries are nonlinear)

# Taxonomy of supervised NL

$$\arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{(\mathbf{x}^i, y^i) \in \mathcal{D}} E(f_{\mathbf{w}}(\mathbf{x}^i), y^i) + R(\dots)$$

$\mathcal{L}(\mathbf{w})$

loss function

set of parameters to estimate " $\mathbf{w}$ "

fidelity to data term

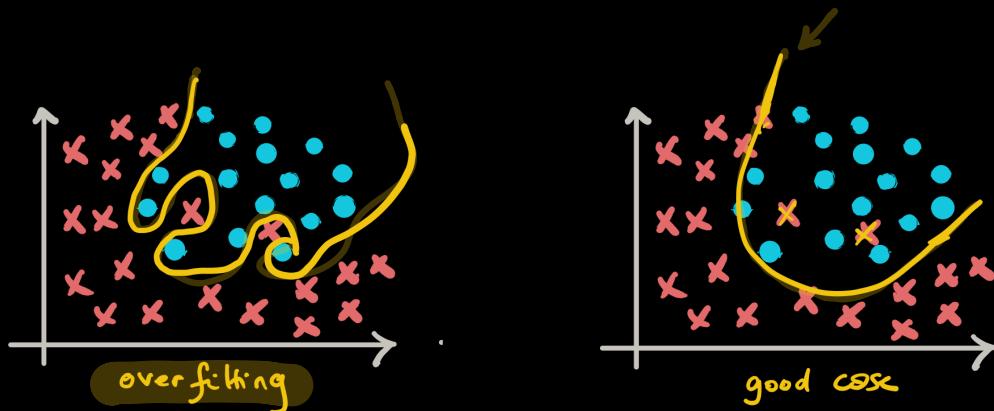
predicted but put by the learner  $f_{\mathbf{w}}$  for input  $\mathbf{x}^i$

error function

regularization term

controls the complexity of the model/learner + avoid overfitting

$d(\hat{y}^i = \tilde{y}^i, y^i)$



I ❤️  $M_x$  [2]

$$\arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{(\mathbf{x}^i, y^i) \in \mathcal{D}} E(f_{\mathbf{w}}(\mathbf{x}^i), y^i) + R(\dots)$$

supervised learning energy cost (loss function)

$f$

$E$

$R$

$\mathbf{w}$

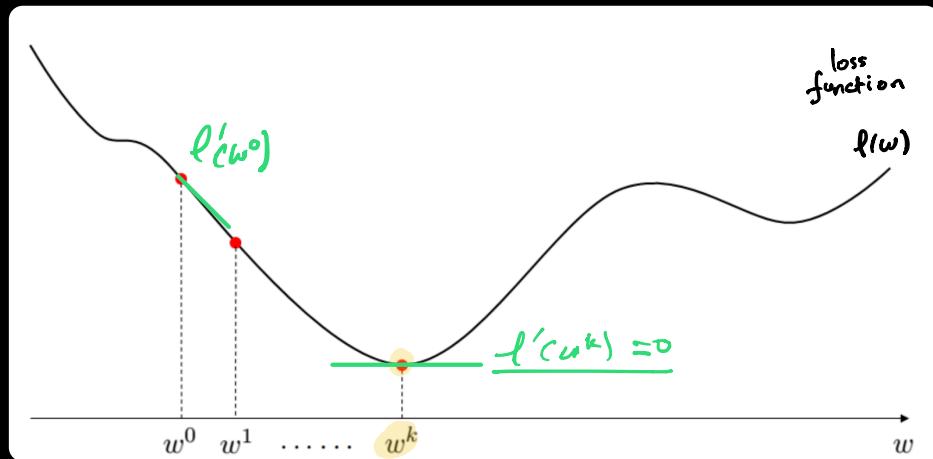
the mapping function to learn from  $\mathbf{x}^i \rightarrow y^i$

the error function between the predicted target by  $f$  and the ground truth observation  $y^i$   
regularization term to avoid overfitting and control model complexity

optimization parameters (weight vector)

set of parameters that minimize the loss functional  $\mathcal{L}(\mathbf{w})$

→ Goal: find the optimal parameters  $w$  for 2 loss function  $\ell$ .

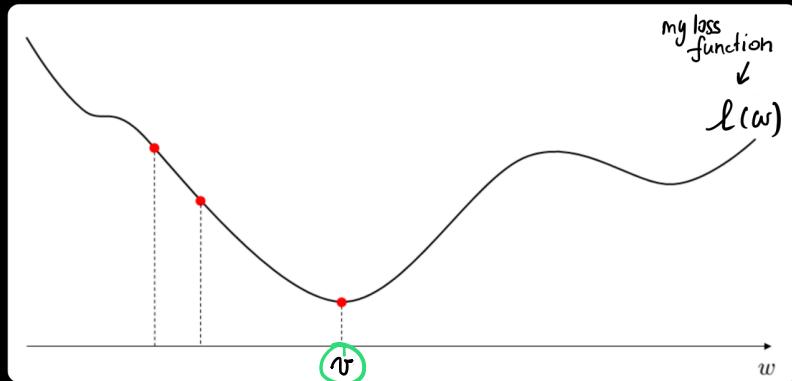


① what are we looking for?

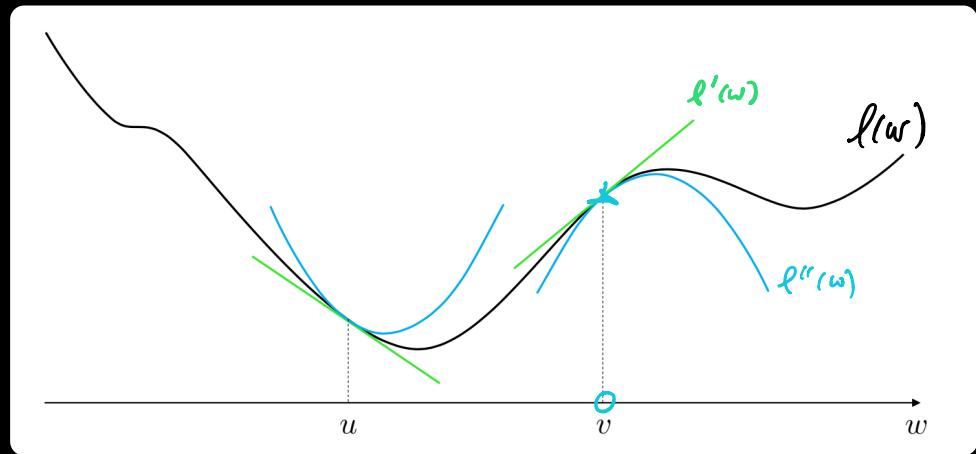
② How to find what we are looking for?

*(hint): "check the local geometry of  $\ell(w)$ "*

How calculus can help us read  
the geometry of 2 function at  
different points?



- 1-dimensional loss function  $\ell(w)$ ,  $w \in \mathbb{R}$ ,  $\ell(w) \in \mathbb{R}$
- To capture the geometry of  $\ell(w)$  at a point  $w^*$ , we can linearly or nonlinearly approximate it near this point.

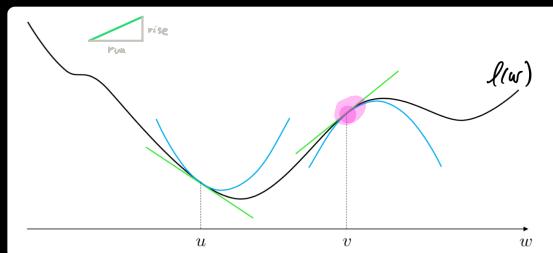


1D

univariate function  $l$  in  $w$ ,  $w \in \mathbb{R}$

**First-order Taylor approximation**

"linear" approximation



• linear approximation of  $l(w)$  at two points  $\begin{cases} w=u \\ w=v \end{cases}$

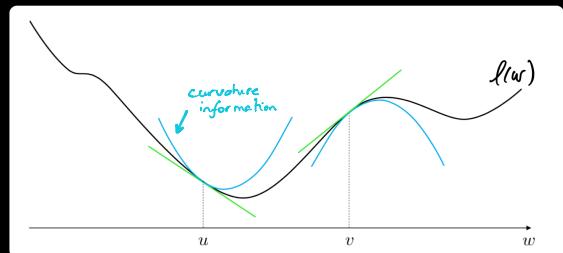
$$h(w) = l(v) + \frac{l'(v)}{\text{rise}}(w-v)$$

$$\begin{cases} h(v) = l(v) \\ h'(v) = l'(v) \end{cases}$$

→ this linear approximation holds well near  $v$  because the derivative contains the slope information.

**Second-order Taylor approximation**

{ "quadratic" approximation  
non-linear



• A second-order Taylor approximation function  $h(w)$  at points  $v$  located near  $w$  contains both first and second derivatives of the target function  $l$  to approximate:

$$h(w) = l(v) + \frac{l'(v)}{\uparrow}(w-v) + \frac{1}{2} \underbrace{l''(v)}_{(w^2 - 2wv + v^2)}(w-v)^2$$

for  $w=v$

$$h(v) = l(v)$$

$$h'(v) = l'(v)$$

$$h''(v) = l''(v)$$

$\left\{ \begin{array}{l} \text{"locally" } \approx \text{ around } w=v \\ \frac{h(w)}{h(v)} \text{ approximates } l(w) \end{array} \right.$

$$\frac{d h(w)}{d w} = h'(w) = l'(v) + \frac{3}{2} \frac{l''(v) w}{2} - \frac{2}{2} \frac{l''(v) v}{2} \Leftarrow$$

$$\Rightarrow h'(w=v) = l'(v) + 0 = l'(v)$$

(ND)

multivariate loss function  $\ell$ ,  
 $w \in \mathbb{R}^{N \times 1}$ ,  $w = [w_1, w_2 \dots w_N]^T$

$h(w) = \ell(v) + \nabla \ell(v)^T (w - v)$

$\nabla \ell(v) = \begin{bmatrix} \frac{\partial}{\partial w_1} \ell(v) \\ \vdots \\ \frac{\partial}{\partial w_N} \ell(v) \end{bmatrix} \in \mathbb{R}^N$

gradient of partial derivatives

$\ell \left( \begin{bmatrix} w \\ w_1 \\ w_2 \end{bmatrix} \right) = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + \begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = 2w_1 + w_2 + w_1^2 + w_2^2$

$\nabla \ell(w) = \begin{bmatrix} \frac{\partial \ell}{\partial w_1} \\ \frac{\partial \ell}{\partial w_2} \end{bmatrix} = \begin{bmatrix} 2+0+2w_1+0 \\ 0+1+0+2w_2 \end{bmatrix} = \begin{bmatrix} 2+2w_1 \\ 1+2w_2 \end{bmatrix}$

$\nabla \ell([1]) = \begin{bmatrix} 4 \\ 3 \end{bmatrix}$

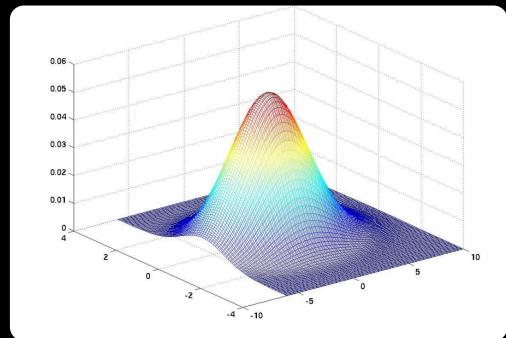
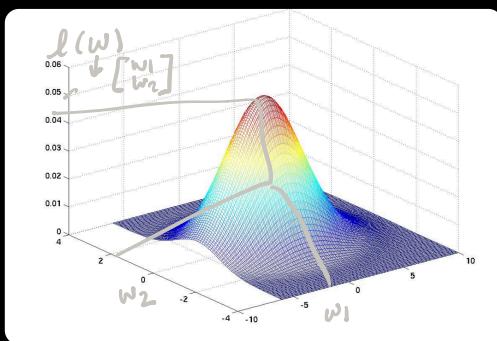
$h(w) = \ell(v) + \nabla \ell(v)^T (w - v) + \frac{1}{2} (w - v)^T \nabla^2 \ell(v) (w - v)$

$\nabla^2 \ell(v) = \begin{bmatrix} \frac{\partial^2}{\partial w_1 \partial w_1} \ell(v) & \cdots & \frac{\partial^2}{\partial w_1 \partial w_N} \ell(v) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial w_N \partial w_1} \ell(v) & \cdots & \frac{\partial^2}{\partial w_N \partial w_N} \ell(v) \end{bmatrix} \in \mathbb{R}^{N \times N}$

$H = \begin{bmatrix} \vdots \\ \cdots \boxed{\frac{\partial^2 \ell}{\partial w_i \partial w_j}} \end{bmatrix}$

\* Special case  $w \in \mathbb{R}^N$ ,  $N=1$

$$\ell''(v) = ?$$



in 2D :  $v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 5 \\ 2 \end{bmatrix} \leftarrow \text{fixed point (evaluate at } v\text{)}$   
 $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \leftarrow \text{our variable}$

$$h(w) = \ell\left(\begin{bmatrix} 5 \\ 2 \end{bmatrix}\right) + \left[\frac{\partial \ell(v)}{\partial w_1} \frac{\partial \ell(v)}{\partial w_2}\right] \times \begin{bmatrix} w_1 - 5 \\ w_2 - 2 \end{bmatrix}$$

$$h(w) = \underbrace{\ell(v)}_{\in \mathbb{R}} + \underbrace{\nabla \ell(v)^T}_{\in \mathbb{R}} \underbrace{(w-v)}_{\in \mathbb{R}}$$

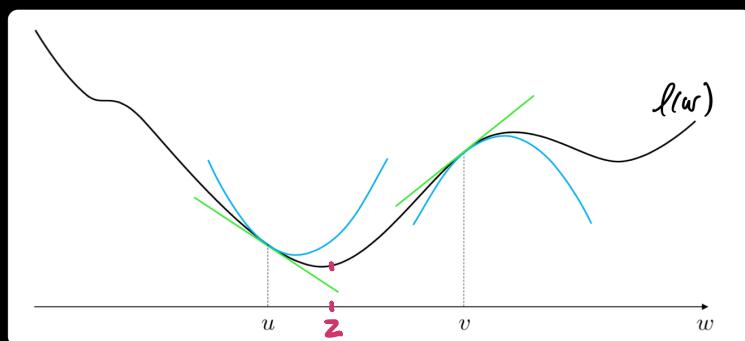
in 2D :  $v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 5 \\ 2 \end{bmatrix} \leftarrow \text{fixed point (evaluate at } v\text{)}$   
 $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \leftarrow \text{our variable}$

$$h(w) = \ell\left(\begin{bmatrix} \quad \\ \quad \end{bmatrix}\right) + \left[ \quad \right] \times \left[ \quad \right]$$

$$+ \left[ \quad \right] \left[ \frac{\partial^2 \ell(v)}{\partial w_1 \partial w_2} \frac{\partial^2 \ell(v)}{\partial w_1 \partial w_2} \right] \times \left[ \quad \right]$$

$$h(w) = \ell(v) + \nabla \ell(v)^T (w-v) + \frac{1}{2} (w-v)^T \nabla^2 \ell(v) (w-v)$$

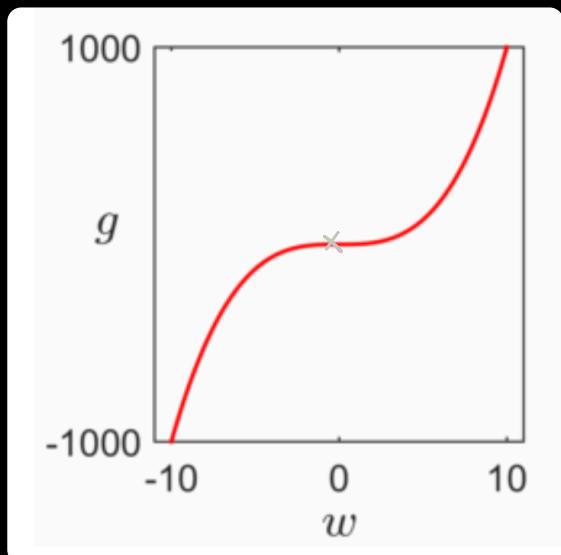
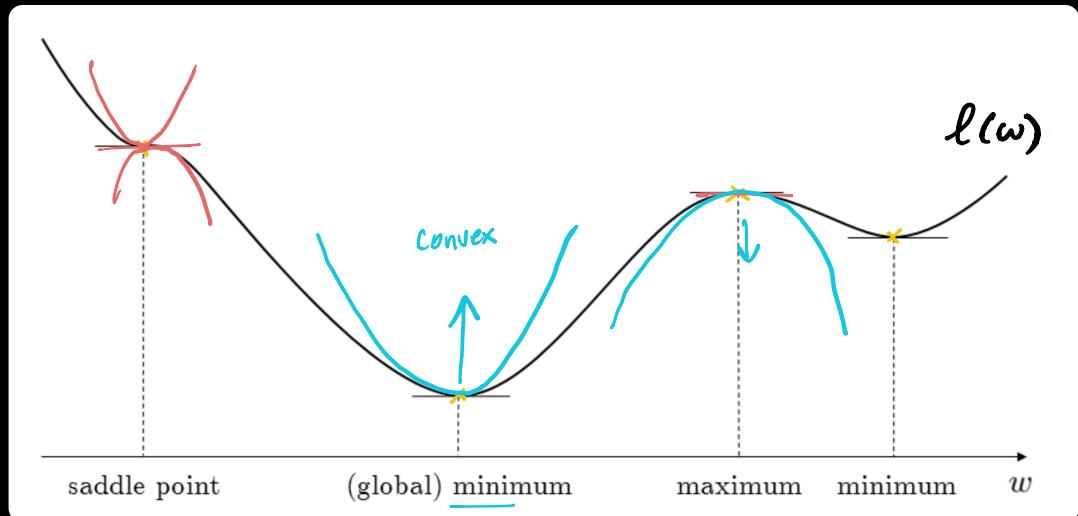
## First order condition of optimality



- Minimum values of  $\ell(w)$  are located at 'valley floors' where the line (or hyperplane in high dimension) tangent to the function  $\ell$  is flat ( $\text{slope} = 0$ ).

$$\begin{cases} N=1 & \Rightarrow \ell'(z) = 0 \\ N>1 & \Rightarrow \nabla \ell(z) = 0_{N \times 1} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \end{cases}$$

$$\nabla \ell(\mathbf{z}) = \begin{bmatrix} \frac{\partial \ell}{\partial w_1}(z) \\ \frac{\partial \ell}{\partial w_2}(z) \\ \vdots \\ \frac{\partial \ell}{\partial w_N}(z) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow \text{first order optimality}$$



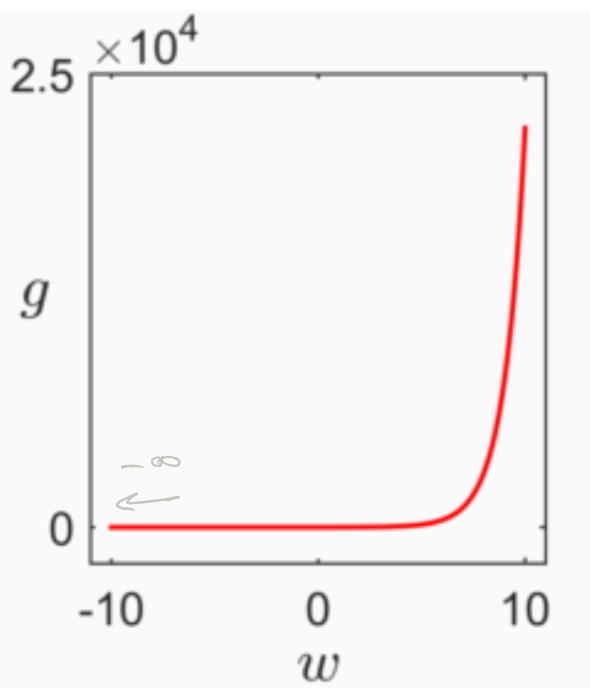
$$g(w) = w^3$$

$$g'(w) = 3w^2$$

$$g'(w) = 0$$

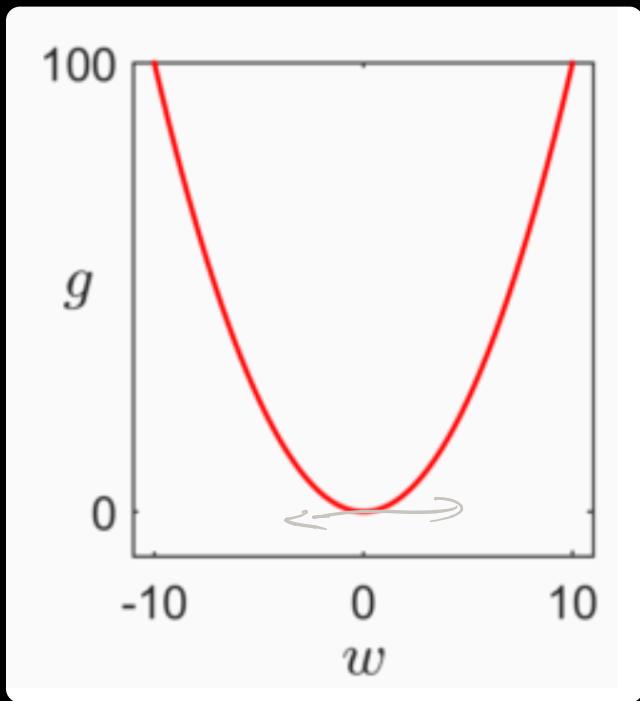
$$w^2 = 0$$

$$\boxed{\Rightarrow w = 0}$$



$$g(\omega) = e^\omega$$

$$\underline{g'(\omega) = e^\omega}$$



$$g(\omega) = \omega^2$$

$$\underline{g'(\omega) = 2\omega}$$

$$g'(0) = 0$$

$$\underline{\omega = 0}$$

$$l(\omega) = \frac{1}{2} \omega^T Q \omega + r^T \omega + d$$



$$\nabla l(\omega) = Q\omega + r$$

$$\Rightarrow \nabla l(\omega) = 0 \Rightarrow \boxed{Q\omega = -r}$$

$$\left\{ \begin{array}{l} Q \in \mathbb{R}^{N \times N} \\ Q^T = Q \quad (\text{symmetric}) \\ r \in \mathbb{R}^{N \times 1} \\ d \in \mathbb{R} \end{array} \right.$$

rules

$$\left\{ \begin{array}{l} \frac{\partial a^T x}{\partial x} = \frac{\partial x^T a}{\partial x} = a \\ \frac{\partial a^T x a}{\partial a} = 2x^T a \\ [x = x^T] = \text{symmetric} \end{array} \right.$$

## Second order condition of optimality

- Goal  $\Rightarrow$  find a global minimum.
- Ideal function is convex [faces upward  $\cup$ ]  $\Rightarrow$  no maxima  $\Rightarrow$  no saddle points



How is the curvature of this function?



1D

$f''(v) > 0 \Leftrightarrow$  convex  
 $f''(v) < 0 \Leftrightarrow$  concave

$N > 1$

multivariate loss functions

$w \in \mathbb{R}^{N \times 1}$ ,  $\ell(w)$  is convex if :

$\nabla^2 \ell(w)$  has positive eigenvalues.

Eigenvectors & eigenvalues  
of a matrix A

- $u$  is an eigenvector of the matrix A if it satisfies:

$$A u = \lambda u$$

$$(A - \lambda I)u = 0$$

eigenvalue      eigenvector

$$A = \begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix}, u = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$A u = \begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 6+4 \\ 6+2 \end{bmatrix} = 4 \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

Example (1st order of a quadratic loss)

$$\ell(w) = \frac{1}{2} w^T Q w + r^T w + d$$

$\frac{\partial \ell}{\partial w} = Qw + r$

$$\nabla \ell(w) = \frac{1}{2} \underbrace{Qw + r}_{=0} \quad \text{rules: } \left\{ \begin{array}{l} Q \in \mathbb{R}^{N \times N} \\ Q^T = Q \quad (\text{symmetric}) \\ r \in \mathbb{R}^{N \times 1} \\ d \in \mathbb{R} \end{array} \right.$$

$$\frac{\partial a^T x}{\partial x} = \frac{\partial x^T a}{\partial x} = a$$

$$\frac{\partial a^T X a}{\partial a} = 2Xa \quad ([x = x^T] = \text{symmetric})$$

$$\Rightarrow \nabla \ell(w) = 0 \Rightarrow Qw + r = 0$$

Example (2nd order of a quadratic loss)

$$\ell(w) = \frac{1}{2} w^T Q w + r^T w + d$$

$$\rightarrow \nabla \ell(w) = \frac{\partial \ell}{\partial w} = Qw + r$$

$$\rightarrow \nabla^2 \ell(w) = \frac{\partial (Qw + r)}{\partial w} \Rightarrow \frac{\partial A x}{\partial x} = \frac{1}{2} (A^T + A)$$

matrix cookbook

$$= \frac{1}{2} (Q^T + Q)$$

$$\rightarrow \nabla^2 \ell(w) = \frac{1}{2} (Q + Q^T) \in \mathbb{R}^{N \times N}$$

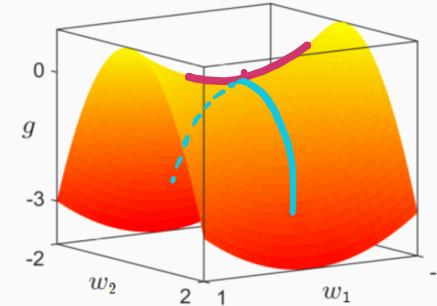
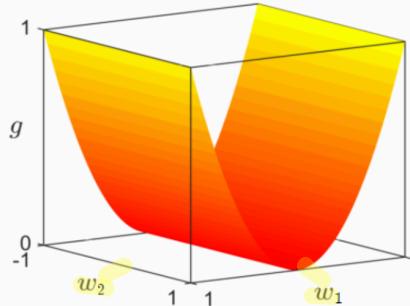
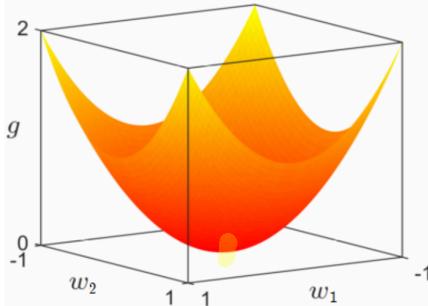
if  $Q$  is symmetric       $Q = Q^T$

$$\nabla^2 \ell(w) = Q$$

$$\ell(\mathbf{w}) = \mathbf{w}^T \mathbf{Q} \mathbf{w}$$

3 examples of convexity in high-dimensional space

$$d=0, \mathbf{r} = \mathbf{0}_{2 \times 1} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



$$\mathbf{Q} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\nabla^2 \ell(\mathbf{w}) = \mathbf{Q}$$

$$\lambda_1 = \lambda_2 = 2$$

$$\mathbf{Q} = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\nabla^2 \ell(\mathbf{w}) = \mathbf{Q}$$

$$\begin{cases} \lambda_1 = 2 \\ \lambda_2 = 0 \end{cases}$$

$$\mathbf{Q} = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

$$\nabla^2 \ell(\mathbf{w}) = \mathbf{Q}$$

$$\begin{cases} \lambda_1 = 2 \\ \lambda_2 = -2 \end{cases}$$

I ❤️ Mx [3]

#### Mathematical notation

#### Definition

$$\arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{(\mathbf{x}^i, y^i) \in \mathcal{D}} E(f_{\mathbf{w}}(\mathbf{x}^i), y^i) + R(\dots)$$

supervised learning energy cost (loss function)

$f$

$E$

$R$

$\mathbf{w}$

the mapping function to learn from  $\mathbf{x}^i \rightarrow y^i$

the error function between the predicted target by  $f$  and the ground truth observation  $y^i$   
regularization term to avoid overfitting and control model complexity  
optimization parameters (weight vector)

set of parameters that minimize the loss function  $\mathcal{L}(\mathbf{w})$

first-order Taylor approximation of the loss function  $l$  at point  $w$  (in 1-dimensional space)  
first derivative of function  $l$  evaluated at point  $w$

second-order Taylor approximation of the loss function  $l$  at point  $w$  (in 1-dimensional space)  
second derivative of function  $l$  evaluated at point  $w$

first-order Taylor approximation of high-dimensional loss function  $l$  at vector point  $\mathbf{w} \in \mathbb{R}^N$   
weight vector to learn

gradient vector of the multivariate loss function  $l$  at location  $\mathbf{w}$

note that  $\nabla l(\mathbf{w})^T \in \mathbb{R}^{1 \times N}$  (row vector)

$$\nabla l(\mathbf{v}) = [\frac{\partial}{\partial w_1} l(\mathbf{v}) \quad \frac{\partial}{\partial w_2} l(\mathbf{v}) \dots \frac{\partial}{\partial w_N} l(\mathbf{v})]^T$$

$h(\mathbf{w}) = l(\mathbf{v}) + \nabla l(\mathbf{w})(\mathbf{w} - \mathbf{v}) + \frac{1}{2}(\mathbf{w} - \mathbf{v})^T \nabla^2 l(\mathbf{w})(\mathbf{w} - \mathbf{v})$  second-order Taylor approximation of high-dimensional loss function  $l$  at vector point  $\mathbf{w} \in \mathbb{R}^N$   
 $\nabla^2 l(\mathbf{w}) \in \mathbb{R}^{N \times N}$  Hessian symmetric matrix of second derivatives of  $l$  along all its dimensions (variables)

$$l'(w) = 0$$

$$\nabla l(\mathbf{w}) = \mathbf{0}_{N \times 1}$$

$l$  is many times differentiable at the vector valued input  $\mathbf{w}$   
stationary point (min, max or saddle) for a 1-dimensional function

stationary point (min, max or saddle) for an  $N$ -dimensional function  
(all elements of the gradient are zero)

$\mathbf{Q}$  is symmetric

matrix cookbook

matrix cookbook

convex function (facing upward)

concave function (facing downward)

matrix cookbook

is the Hessian matrix of  $l$  equal to  $l(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{Q} \mathbf{w} + \mathbf{r}^T \mathbf{w} + d$   
 $d \in \mathbb{R}$  and  $\mathbf{r} \in \mathbb{R}^{N \times 1}$

is the gradient vector of  $l$  equal to  $l(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{Q} \mathbf{w} + \mathbf{r}^T \mathbf{w} + d$

$$\mathbf{Q} = \mathbf{Q}^T$$

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}$$

$$\frac{\partial \mathbf{a}^T \mathbf{x} \mathbf{a}}{\partial \mathbf{a}} = 2 \mathbf{x} \mathbf{a} (\mathbf{x}^T \mathbf{x})$$

$$l''(w) > 0$$

$$l''(w) < 0$$

$$\frac{\partial \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \frac{1}{2} (\mathbf{A}^T + \mathbf{A})$$

$$\nabla^2 l(\mathbf{w}) = \frac{1}{2} (\mathbf{Q}^T + \mathbf{Q}) \in \mathbb{R}^{N \times N}$$

$$\nabla^2 l(\mathbf{w}) = \mathbf{Q} \mathbf{w} + \mathbf{r} \in \mathbb{R}^{N \times 1}$$

Basic technique 1 for finding  $\min_w \ell(w)$ :  
gradient descent

$$w^* = \underset{w}{\text{minimize}} \quad \underline{\ell(w)}$$

minimize  $\ell$  over all input values  $w$



solution = "optimal  $w$ "



$w^*$

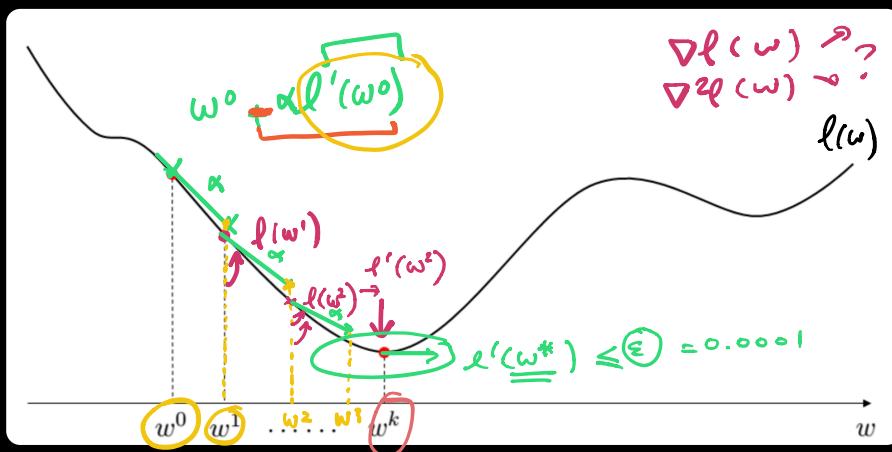


Solve  $N$  equations

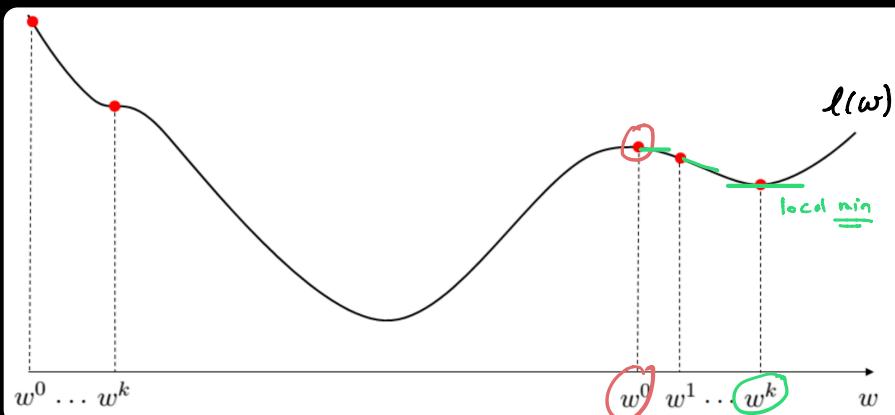
$$\boxed{\nabla \ell(w) = 0_{N \times 1}}$$

choice of  $w_0$

non-convex

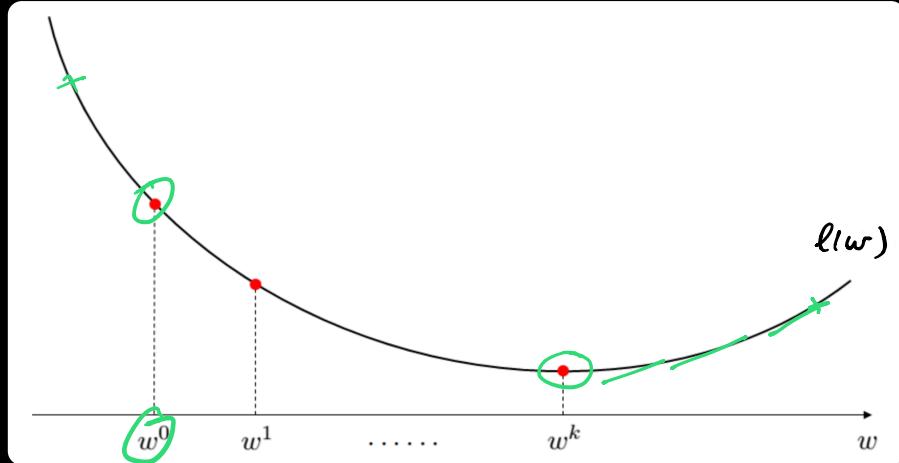


non-convex



influence of  
 $w^0$  ??

convex or  
non convex?

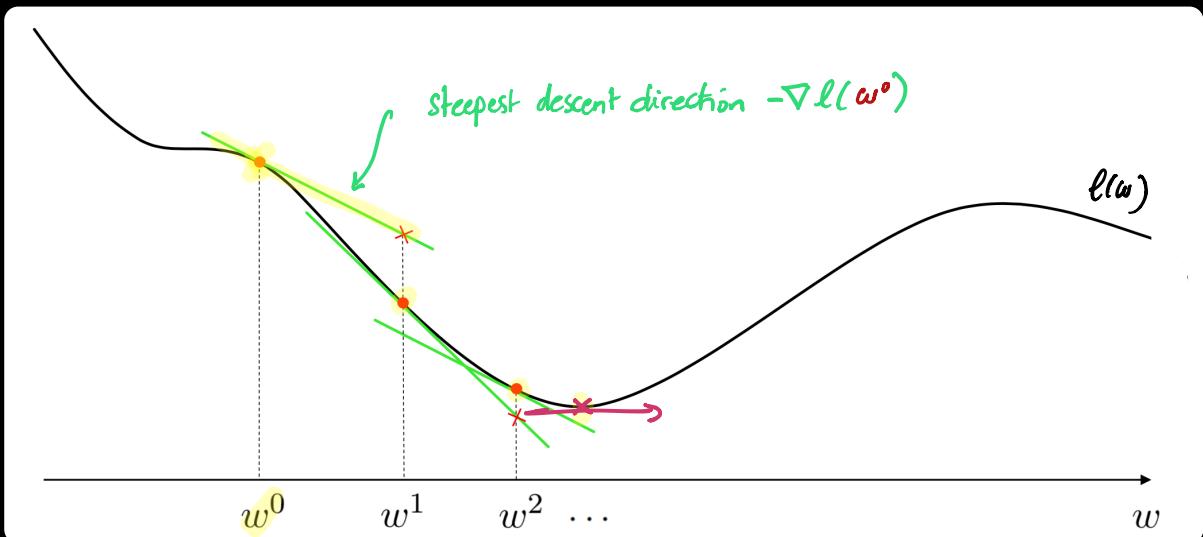


- ① Start the minimization process from some initial point  $\mathbf{w}^0$ .
- ② Take iterative steps denoted by  $\mathbf{w}^1, \mathbf{w}^2, \dots$ , going “downhill” towards a stationary point of  $\mathbf{l}$ .
- ③ Repeat step ② until the sequence of points converges to a stationary point of  $\mathbf{l}$ .

### STOPPING CONDITION

- ① When a pre-specified number of iterations are complete.
- ② When the gradient is small enough, i.e.,  $\|\nabla \mathbf{l}(\mathbf{w}^k)\|_2 < \epsilon$  for some small  $\epsilon > 0$ .

## Gradient descent



- ① Travel in the downward direction of a linear approximation
- ② hop back onto the function
- ③ repeat in order to find a stationary point of  $\ell$ .

$$* \underline{h}(w) = \underline{\underline{l}}(w^0) + \nabla \underline{l}(w^0)^T (w - w^0)$$

→ take a first step of size  $\alpha_1$  downwards:

$$\Rightarrow \underline{\circlearrowleft} \underline{\underline{w}}^1 = \underline{\underline{w}}^0 - \underline{\underline{\alpha}}_1 \underline{\underline{\nabla}} \underline{\underline{l}}(\underline{\underline{w}}^0)$$

learning rate / step length

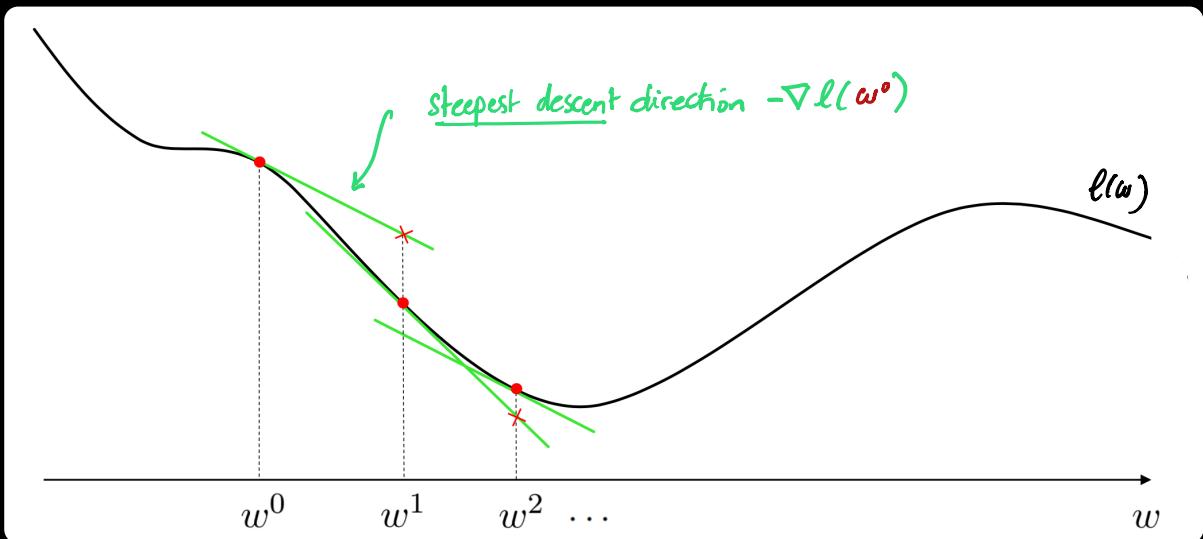
$$\underline{\underline{w}}^2 = \underline{\underline{w}}^1 - \underline{\underline{\alpha}}_2 \underline{\underline{\nabla}} \underline{\underline{l}}(\underline{\underline{w}}^1)$$

$$\Rightarrow \underline{\circlearrowleft} \underline{\underline{w}}^k = \underline{\underline{w}}^{k-1} - \underline{\underline{\alpha}}_k \underline{\underline{\nabla}} \underline{\underline{l}}(\underline{\underline{w}}^{k-1})$$

$\downarrow$   
 $w^*$   $k^{\text{th}}$  gradient descent

generally  $\Rightarrow \underline{\underline{\alpha}}_k = \underline{\underline{\alpha}}$

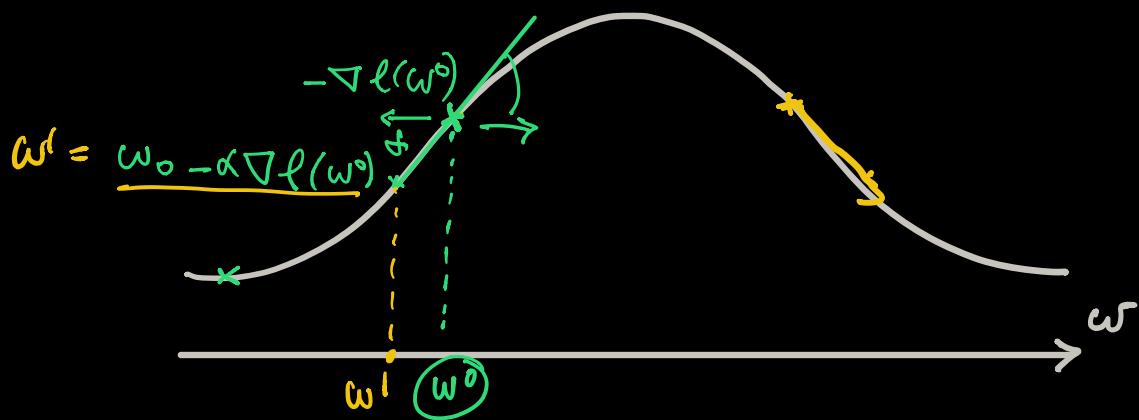
## Gradient descent



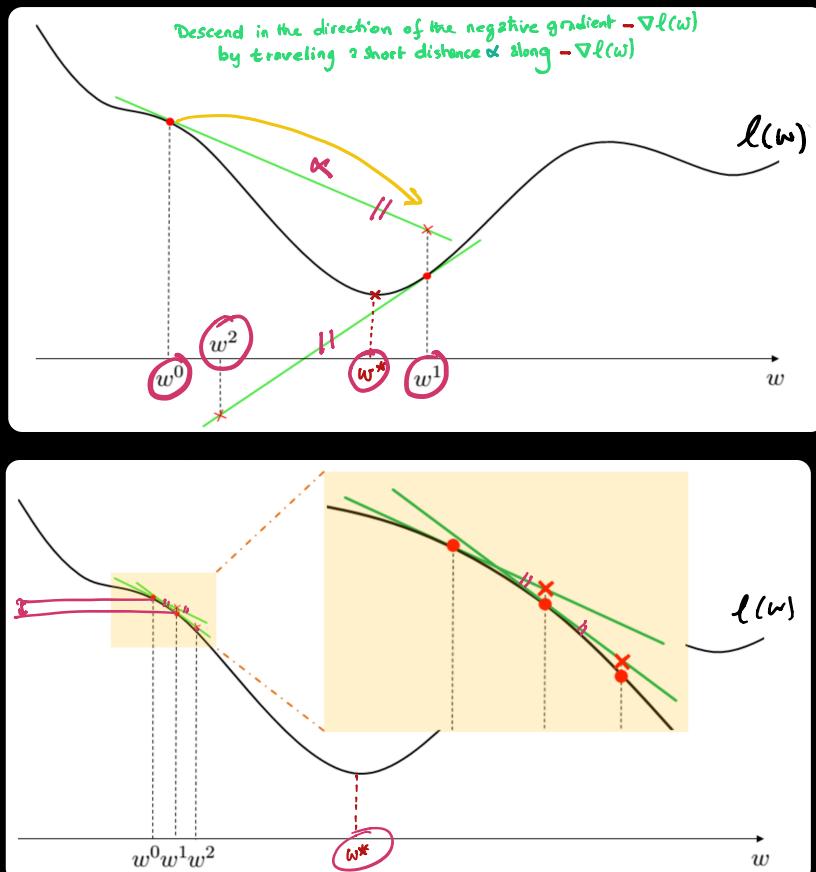
- ① Travel in the downward direction of a linear approximation
- ② hop back onto the function
- ③ repeat in order to find a stationary point of  $l$ .

$$h(w) = l(w^0) + \nabla l(w^0)^T (w - w^0)$$

“How to find The direction with the steepest descent ? ”



How to choose  $\alpha$  ?



### Algorithm 2.1 Gradient descent (with fixed step length)

**Input:** differentiable function  $\ell$ , fixed step length  $\alpha$ , and initial point  $\mathbf{w}^0$

$k = 1$

$$\nabla \ell(\mathbf{w}^k) \leq \varepsilon$$

Repeat until **stopping condition** is met:

$$\mathbf{w}^k = \mathbf{w}^{k-1} - \alpha \nabla \ell(\mathbf{w}^{k-1})$$

$$k \leftarrow k + 1$$

[Matlab / Python demo]