

# Gradient Descent in action

Class starts @ 9:03 PM

- Biz-case
- deadline
- Biz-case
- discussion

2 Weeks



=====

Agenda:

Math:  $f(x)$

① Gradient Descent code + experiments

{ ② Linear Regression as optimization

→ ML-1  
4-cl

③ Revisit the classification problem

+ Code + Text

RAM Disk

```
0s
import random
import math
import numpy
```

```
[ ] #define f(x)
def f(x):
    return math.pow(x, 2)+2*x+10

# df_dx; we could have used the limits concepts from prev class too.
def derF(x):
    return 2*x+2
```

$$\frac{dx^n}{dx} = nx^{n-1}$$

Math:

One-variable

$$f(x) = x^2 + 2x + 10$$

$$\min_x f(x)$$

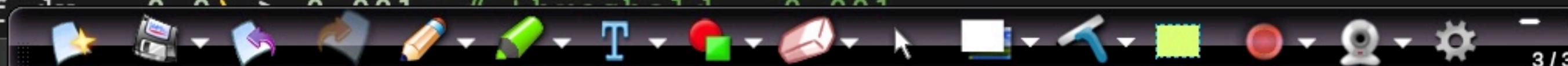
$$\begin{aligned}\frac{df(x)}{dx} &= 2x + 2 \cdot 1 + 0 \\ &= 2x + 2\end{aligned}$$

```
def gradientDescentUnivariate():

    x0 = random.random()
    i=0 # iterations
    eta = 0.1 # learning rate

    xi = x0
    df_dx = derF(xi) # df_dx at x0

    while abs(df_dx) > 0.001:
```



+ Code + Text

```
return 2*x+2
```

```
[ ] def gradientDescentUnivariate():
```

```
x0 = random.random()
```

```
i=0 # iterations
```

```
eta = 0.1 # learning rate
```

```
xi = x0
```

```
df_dx = derF(xi) # df_dx at x0
```

```
while abs(df_dx - 0.0) > 0.001: # threshold = 0.001
```

```
{ print(i)
  print(xi)
  print(df_dx)
  print("*****") }
```

```
{ df_dx = derF(xi) # compute the derivative
```

```
xi = xi - eta*df_dx # update equation
```

```
i = i+1
```

```
return xi
```

$$\min_x f(x)$$

$x_0$ : randomly

$x_1$

$x_2$

$x_3$

⋮

$$x_{t+1} = x_t - \eta \frac{df}{dx}$$

[+ Code](#)[+ Text](#)

return 2\*x+2

{x} [ ] def gradientDescentUnivariate():

{ x0 = random.random()  
i=0 # iterations  
eta = 0.1 # learning ratexi = x0  
df\_dx = derF(xi) # df\_dx at x0

{ while abs(df\_dx - 0.0) &gt; 0.001: # threshold = 0.001

    print(i)  
    print(xi)  
    print(df\_dx)  
    print("\*\*\*\*\*")    df\_dx = derF(xi) # compute the derivative  
    xi = xi - eta\*df\_dx # update equation  
    i = i+1

return xi

✓ RAM Disk



+ Code + Text

```
# iterations
eta = 0.1 # learning rate

xi = x0
df_dx = derF(xi) # df_dx at x0

while abs(df_dx - 0.0) > 0.001: # threshold = 0.001
    print(i)
    print(xi)
    print(df_dx)
    print("*****")
    df_dx = derF(xi) # compute the derivative
    xi = xi - eta*df_dx # update equation
    i = i+1

return xi
```

$$f_{\min} = \min_x f(x)$$

$$x_{\min} = x^* = \tilde{\arg \min}_x f(x)$$

```
x_star = gradientDescentUnivariate()
print(x_star)
print(f(x_star))
```

```
0
0.7326180598960425
3.4652361197920847
```

+ Code + Text

```
# iterations
eta = 0.1 # learning rate

xi = x0
df_dx = derF(xi) # df_dx at x0

while abs(df_dx - 0.0) > 0.001: # threshold = 0.001
    print(i)
    print(xi)
    print(df_dx)
    print("*****")
```

```
df_dx = derF(xi) # compute the derivative
xi = xi - eta*df_dx # update equation
i = i+1
```

```
return xi
```

```
x_star = gradientDescentUnivariate()
print(x_star)
print(f(x_star))
```

```
0
0.7326180598960425
3.4652361197920847
```



GradientDescent.ipynb - Colab x Quotient rule - Wikipedia x Chain rule - Wikipedia x +

colab.research.google.com/drive/1WUN4yW7\_uFW2kW8cc03zwKc\_Rj1E25f8#scrollTo=wzrMvzPkvlxl

+ Code + Text

RAM Disk

0s

0

0.4626230734029989 ↗

2.925246146805998

{x}

1 ✓

0.170098458722391 ✓

2.925246146805998

\*\*\*\*\*

2

-0.06392123302208072

2.340196917444798

\*\*\*\*\*

3

-0.2511369864176646

1.8721575339558385

\*\*\*\*\*

4

-0.40090958913413166

1.4977260271646708

\*\*\*\*\*

5

-0.5207276713073054

1.1981808217317367

\*\*\*\*\*

6

-0.6165821370458443

0.9585446573852893

i  
x\_i  
 $\frac{df}{dx} |_{x_i}$

8 / 10



+ Code + Text

0s

0.002897023463917181  
\*\*\*\*\*  
33  
-0.9990729524915465  
0.0023176187711337892  
\*\*\*\*\*  
34  
-0.9992583619932371  
0.0018540950169070314  
\*\*\*\*\*  
35  
-0.9994066895945897  
0.001483276013525714  
\*\*\*\*\*  
36  
-0.9995253516756717  
0.00118662081082066  
\*\*\*\*\*  
-0.9996202813405374  
9.000000144186261

-1 =  $x^*$

$f(x)$

$\min_x f(x) = q$

$\arg \min_x f(x) = -1$

! 26s

19994  
0.30312655233365504  
-2.60625310466731  
\*\*\*\*\*  
19995  
-2.303126552333655 → 2  
2.60625310466731  
\*\*\*\*\*  
19996  
0.30312655233365504 → 3  
-2.60625310466731  
\*\*\*\*\*  
19997  
-2.303126552333655 → 4  
2.60625310466731  
\*\*\*\*\*  
19998  
0.30312655233365504  
-2.60625310466731  
\*\*\*\*\*  
19999  
-2.303126552333655  
2.60625310466731  
\*\*\*\*\*  
20000  
0.30312655233365504  
-2.60625310466731  
\*\*\*\*\*

$$0.00 \leftarrow \rightarrow \gamma = 0.1 \rightarrow 36$$

$$\eta = 0.3 \rightarrow 9$$

$$\gamma = 0.5 \rightarrow 1$$

$$\left\{ \gamma = 1.0 \rightarrow \infty \right. \text{ (oscillation problem)}$$

[ hack:  $\eta_{i+1} = \eta_i + 0.9$

Q

$$f(x,y) = x^2 + y^2 + 2xy + 5x + 3y + 2$$

9:39 AM

Obj: find  $x^*$  &  $y^*$  that  
 $\min_{x,y} f(x,y)$

i

$$\frac{\partial f}{\partial x} = 2x + 2y + 5$$

$$\frac{\partial f}{\partial y} = 2y + 2x + 3$$

Pseudo code:

$x_i = \text{random-number}$

$y_i = \text{random-number}$

$i = 0 ; \eta = 0.1$

✓  $\partial f - \partial x = \text{derfx}(x_i)$

✓  $\partial f - \partial y = \text{derfy}(y_i)$

{ NOT  
BOUNDARY  
cases

$$x_i : \frac{\partial f}{\partial x} \approx 0$$

$$y_i : \frac{\partial f}{\partial y} \approx 0$$

while  $\left( \frac{\text{abs}(\partial f - \partial x)}{\partial f - \partial x} > 0.001 \right) \text{ OR } \left( \frac{\text{abs}(\partial f - \partial y)}{\partial f - \partial y} > 0.001 \right)$

{  
 $\partial f - \partial x = \text{derfx}(x_i)$   
 $\partial f - \partial y = \text{derfy}(y_i)$

$$\left\{ \begin{array}{l} x_i = x_i - \eta (\partial f - \partial x) \\ y_i = y_i - \eta (\partial f - \partial y) \\ i = i + 1 \\ \dots \\ \text{return } (x_i, y_i) \end{array} \right.$$

# Boundary Case

①

saddle point  $\rightarrow$

$(x_i, y_i)$

(last-class)  
perturbation  
test

②

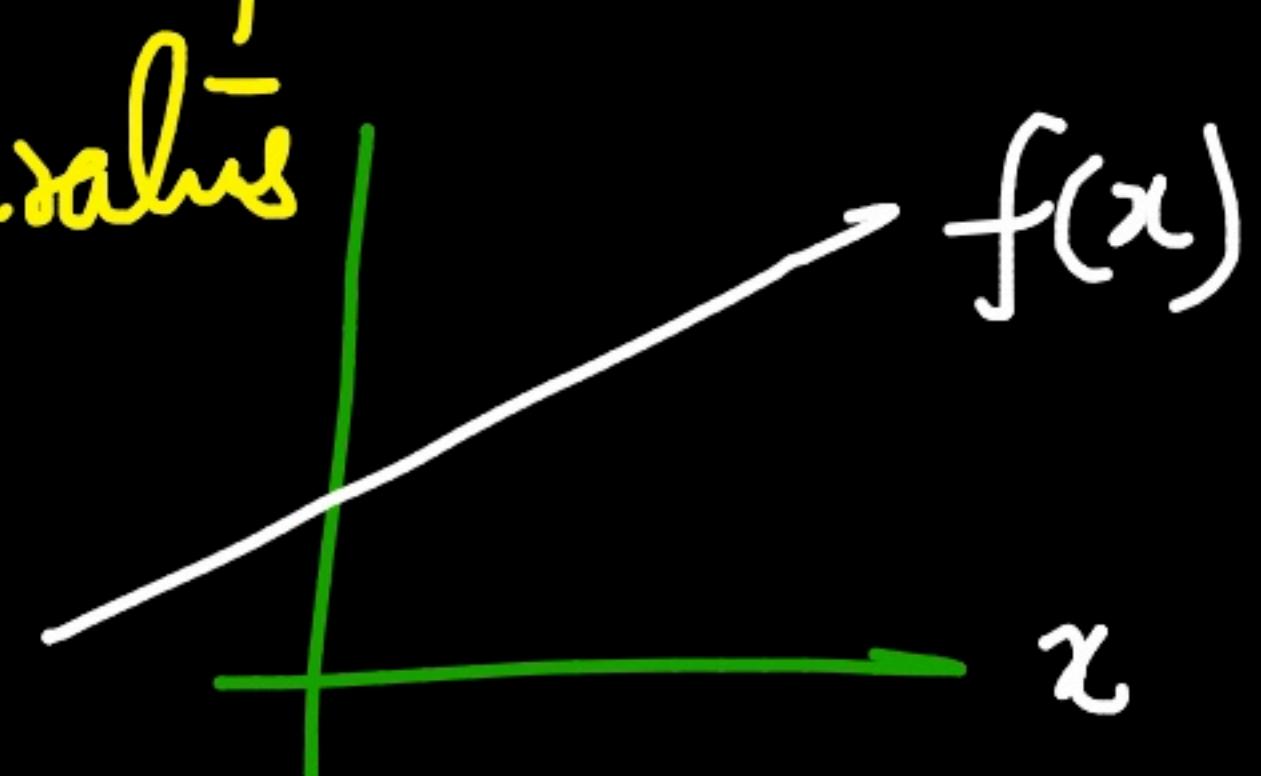
no minima  $\rightarrow$  look iterates

③

oscillation  $\rightarrow x_i, x_{i+1}, x_{i+2}$

④

0.001: threshold  $\rightarrow 1e^{-6} \text{ or } 1e^{-9}$



+ Code + Text

80982  
-6.69756756696005  
0.0012339248835084062  
\*\*\*\*\*

{x}  
80983  
-6.697579906056629  
0.001233909657896161  
\*\*\*\*\*  
80984  
-6.697592245000956  
0.0012338944326596554  
\*\*\*\*\*

80985  
-6.697604583793034  
0.0012338792077988753  
\*\*\*\*\*  
80986  
-6.697616922432867  
0.001233863983313807

$$\frac{d e^x}{dx} = \overset{\sim}{e^x}$$

Artificial  
Threshold

0.001

10<sup>-6</sup>

10<sup>9</sup>

$$\frac{\partial f}{\partial x} \rightarrow 0$$

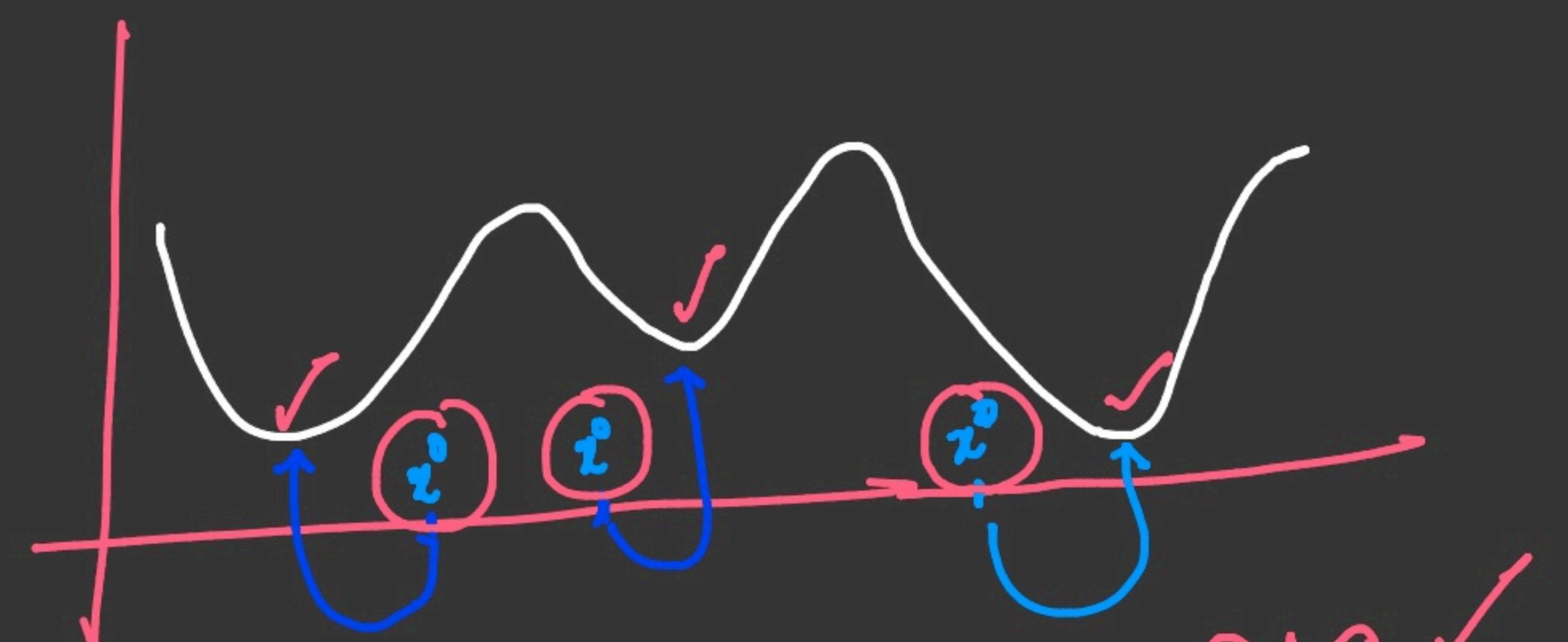
$$\frac{\partial f}{\partial y} \neq 0$$

$$x_{i+1} = x_i - n \left( \frac{\partial f}{\partial x} \right)$$

$$y_{i+1} = y_i - n \left( \frac{\partial f}{\partial y} \right)$$

+ Code + Text

```
*****
2m 99936
-6.907728433272202
0.0010000368467572279
*****
{x}
99937
-6.907738433540662
0.0010000268460702846
*****
99938
-6.907748433709118
0.0010000168455833601
*****
99939
-6.9077584337775715
0.0010000068452964478
*****
-6.907768433746024
0.0009999868453226368
```

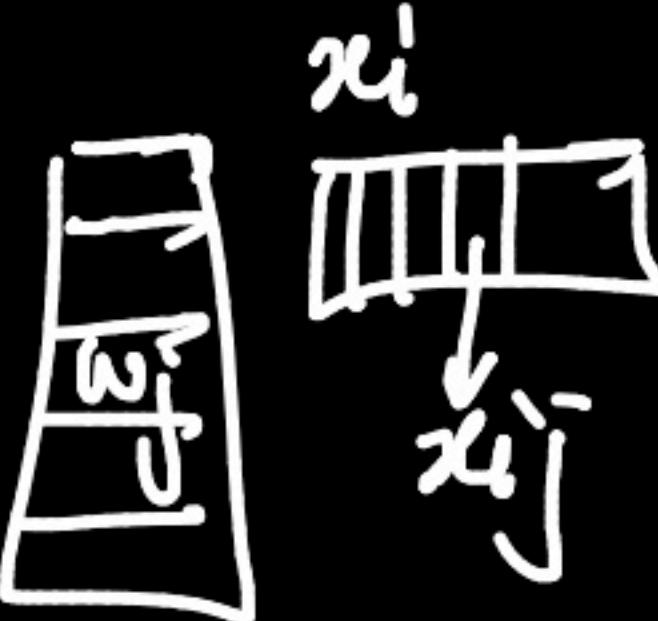


~  
200

✓ one-minima → convex fns

✓ reach any minima

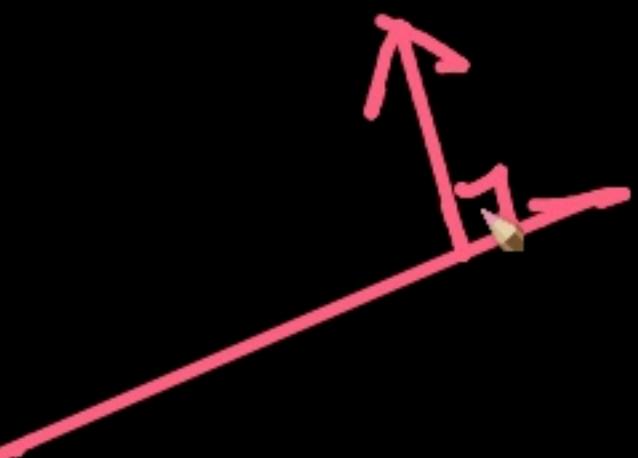
← non-convex fns  
DL



classification

$$\text{Max} \sum_{i=1}^n \frac{\omega^\top x_i + b}{\|\omega\|} \cdot y_i$$

$$\min_{\substack{\max \\ w_j, b}} - \sum_{i=1}^n \left( \frac{\sum_{j=1}^d w_j x_{ij} + b}{\|\omega\|} \right) \cdot y_i \rightarrow L(w_1, w_2, \dots, w_d, b)$$



$d+1$  variables

derivatives:

$$\frac{\partial L}{\partial w_j}$$

for all  $j=1 \dots d$

$$\frac{\partial L}{\partial b}$$

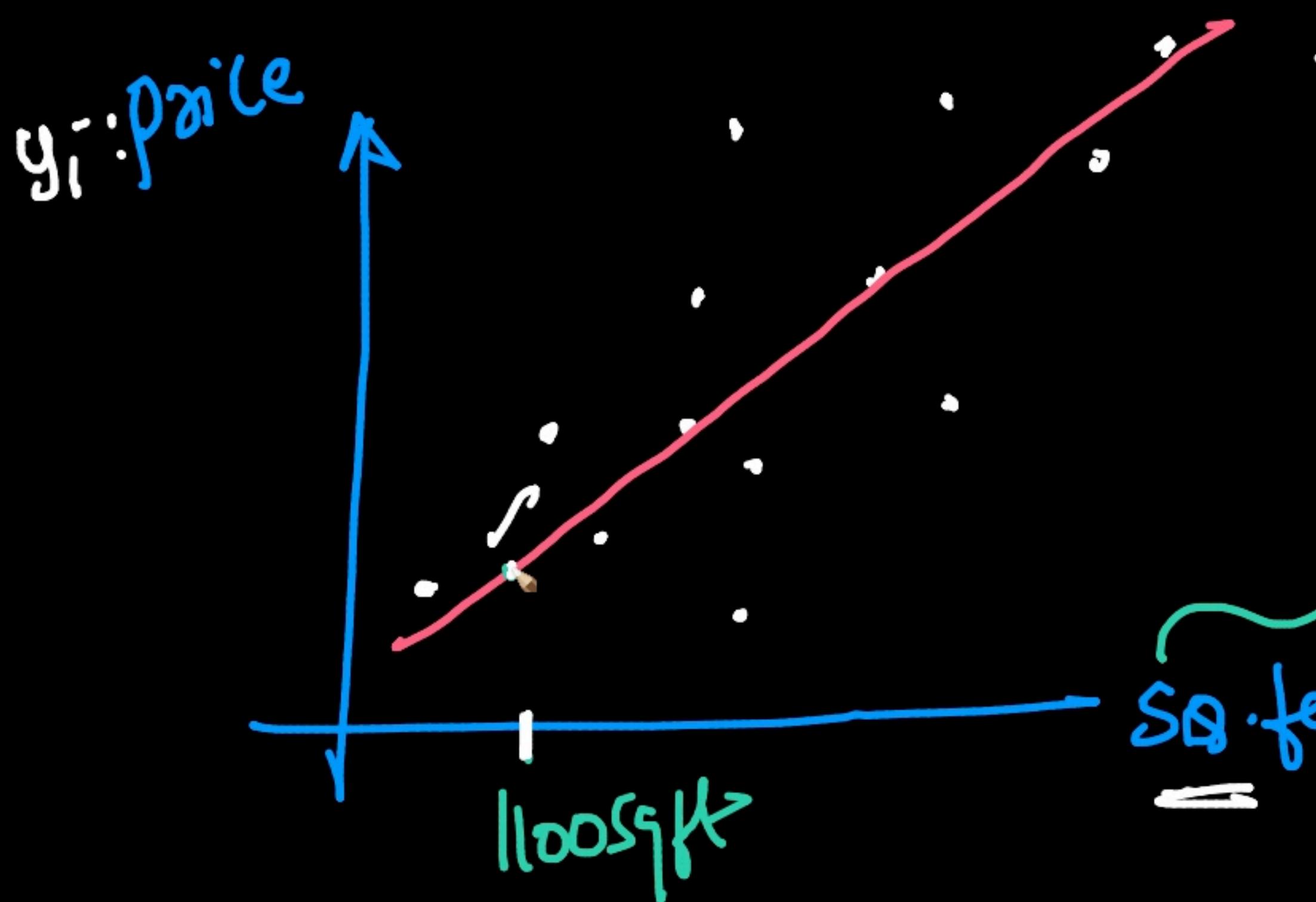
$$\frac{d}{dw_j} \left( \sum_{j=1}^d w_j \cdot x_{ij} + b \right)$$

Cumbersome but  
doable

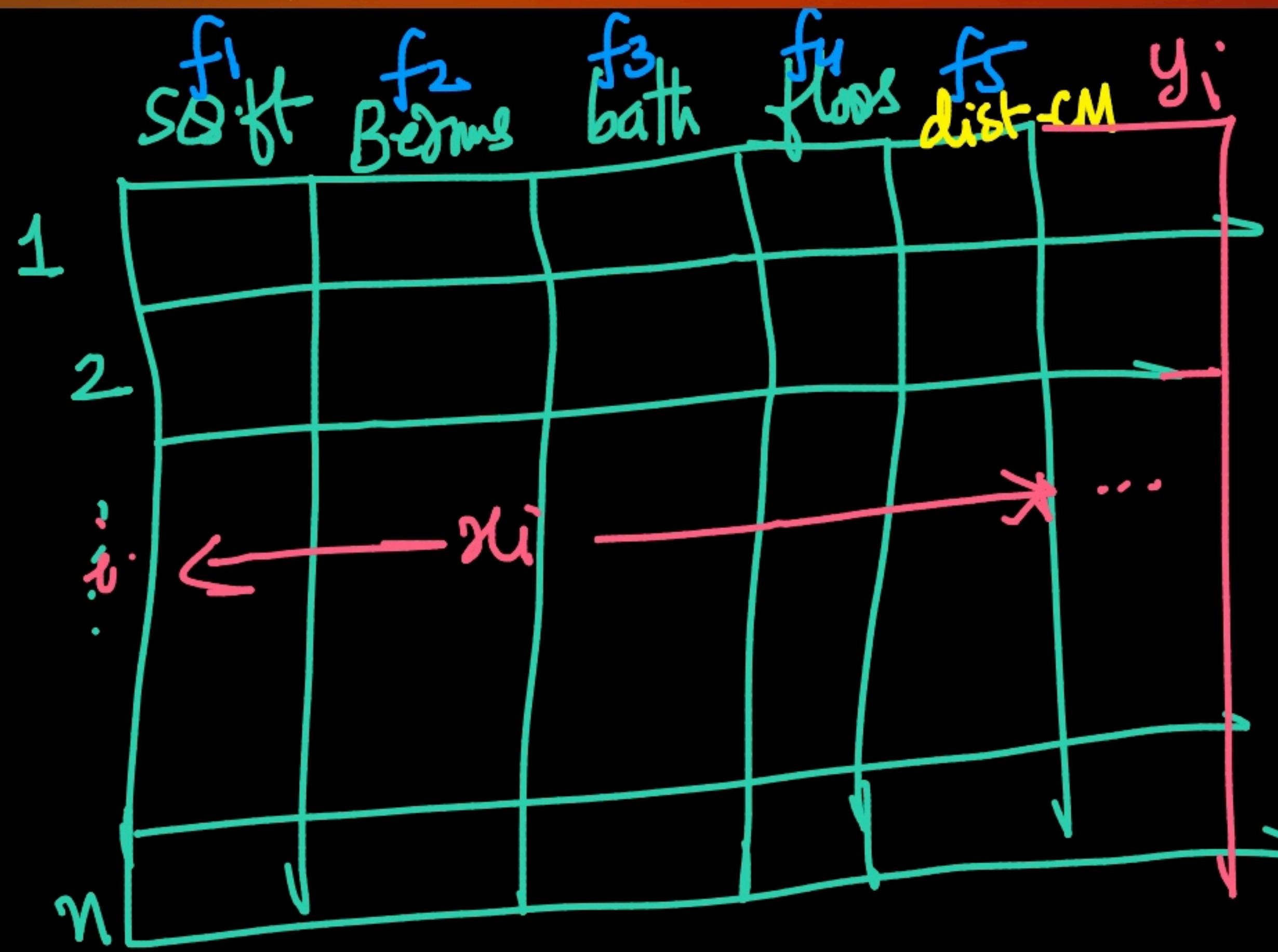
linear regression

$i=1 \text{ to } n$

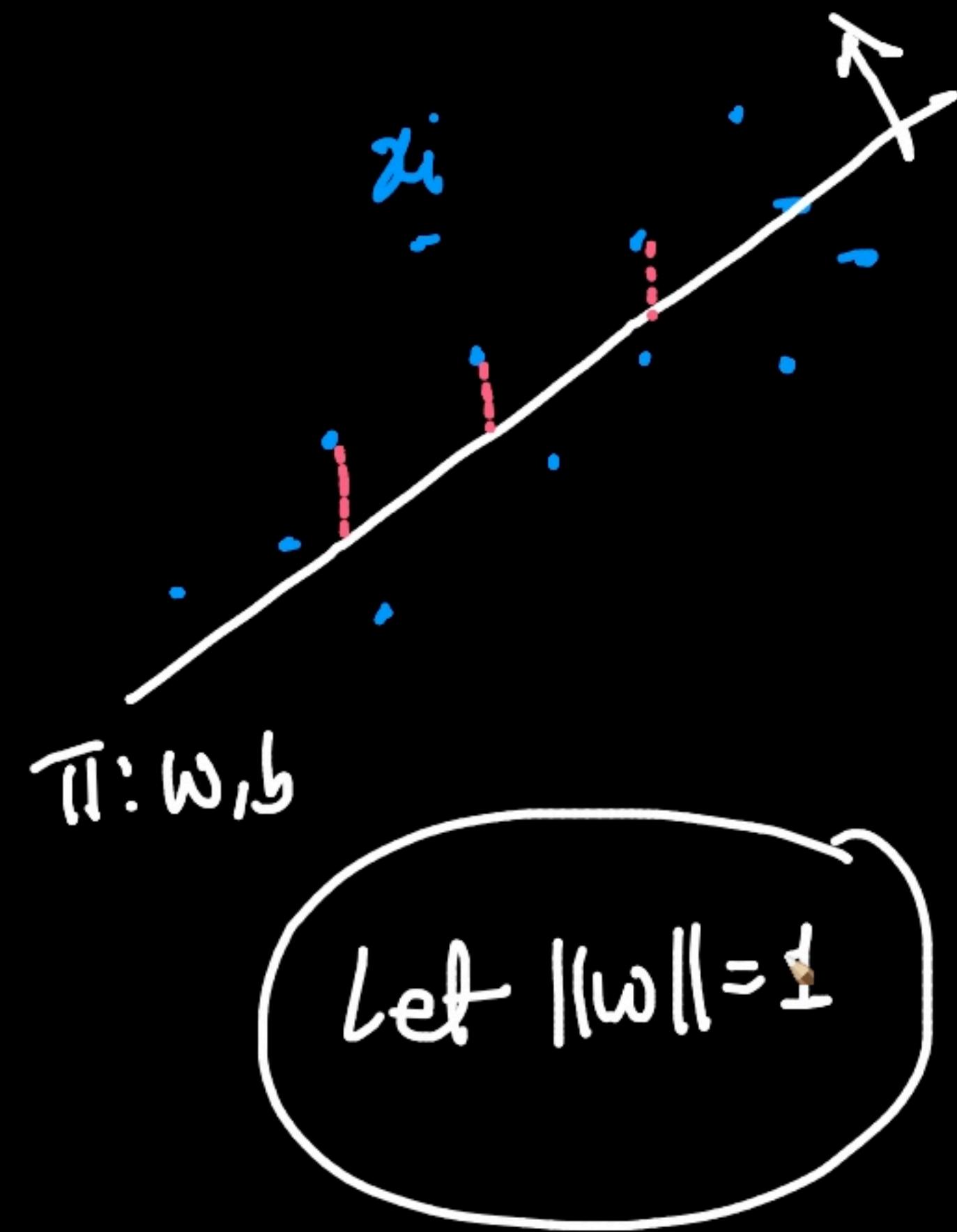
$\{x_i, y_i\}$



obj: line that  
is as close  
as possible  
to given  
data



$$y_i = \vec{w}^T \vec{x}_i + b$$



$$\min_{\underline{w}, b} \sum_{i=1}^n \left( y_i - (\underline{w}^\top \underline{x}_i + b) \right)^2$$

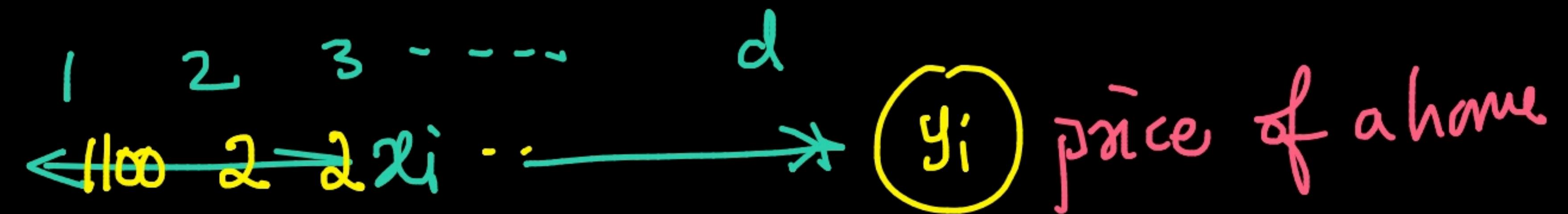
*actual*      *pred*

$L(w, b)$

$\frac{\partial L}{\partial w_j} = \sum_{i=1}^n (-x_{ij}) \cdot 2(y_i - (\underline{w}^\top \underline{x}_i + b))$

$\frac{\partial L}{\partial b} = \sum_{i=1}^n (-1) \cdot 2(y_i - (\underline{w}^\top \underline{x}_i + b))$

$y_i - (\underline{w}_1 x_{i1} + \underline{w}_2 x_{i2} + \dots + \underline{w}_d x_{id} + b)^2$



$$y_i = w_0 x_{i1} + w_1 x_{i2} + w_2 x_{i3} + \dots + w_d x_{id} + b$$

pred - value:

$$\hat{w}^T x_i + b$$

## Recap:

①  $f(x)$ : one variable  $\rightarrow$  code

②  $f(x, y) \rightarrow$  2-var  $\rightarrow$  pseudo-code  
+ boundary cases

③  $L(w_1, w_2, \dots, w_D, b) \rightarrow \checkmark$

④ Linear Regr  $\rightarrow$  optim 2n

ML

$\sum_{i=1}^n$  : Computational expensive

Next class:

$$\min_{w,b} \mathcal{L}(w, b)$$

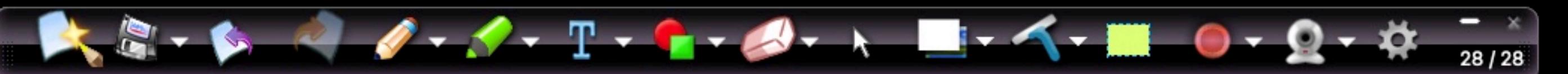
vec  
scalar

s.t.  $\|w\| = 1$

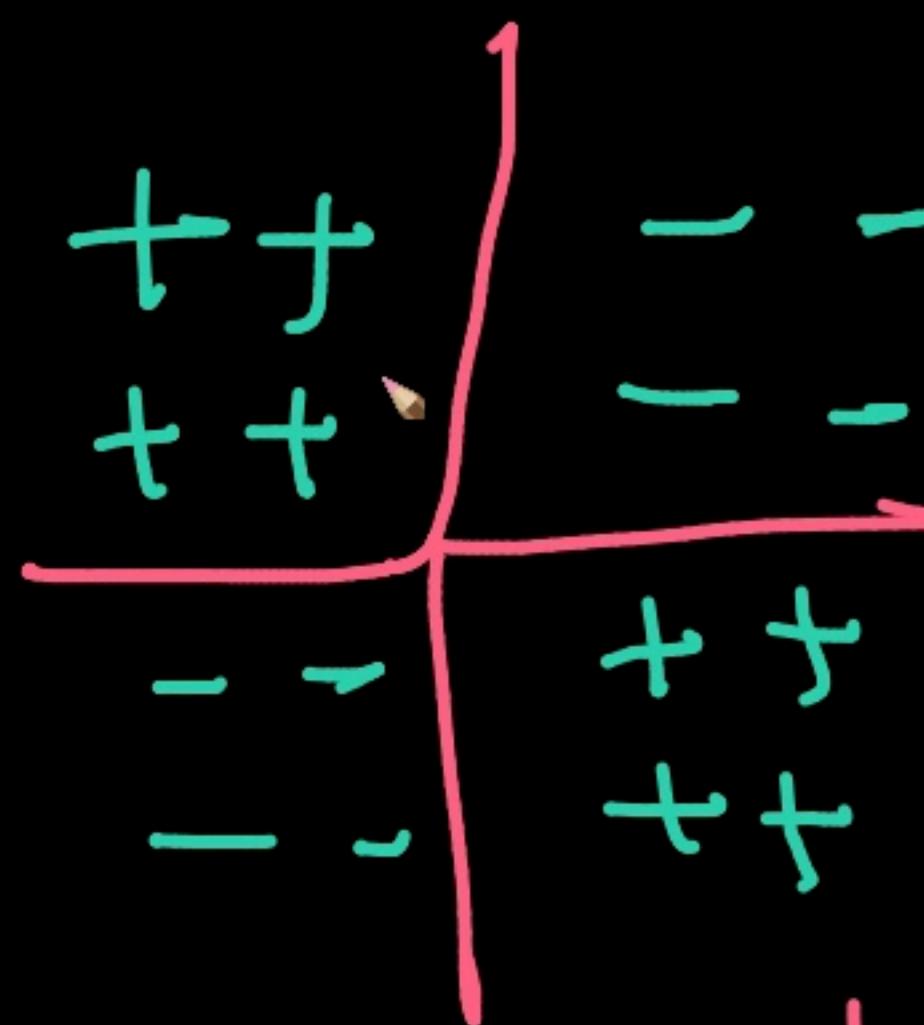
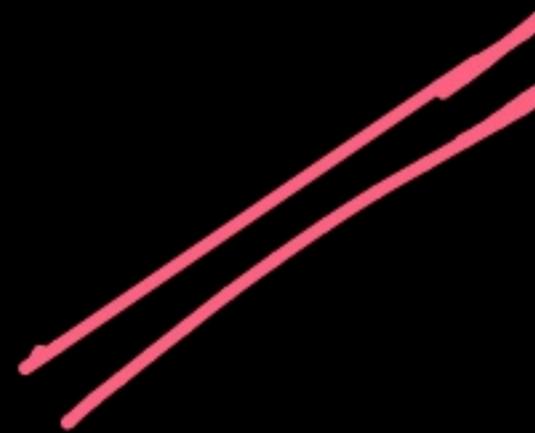
classif

Constrained optimization

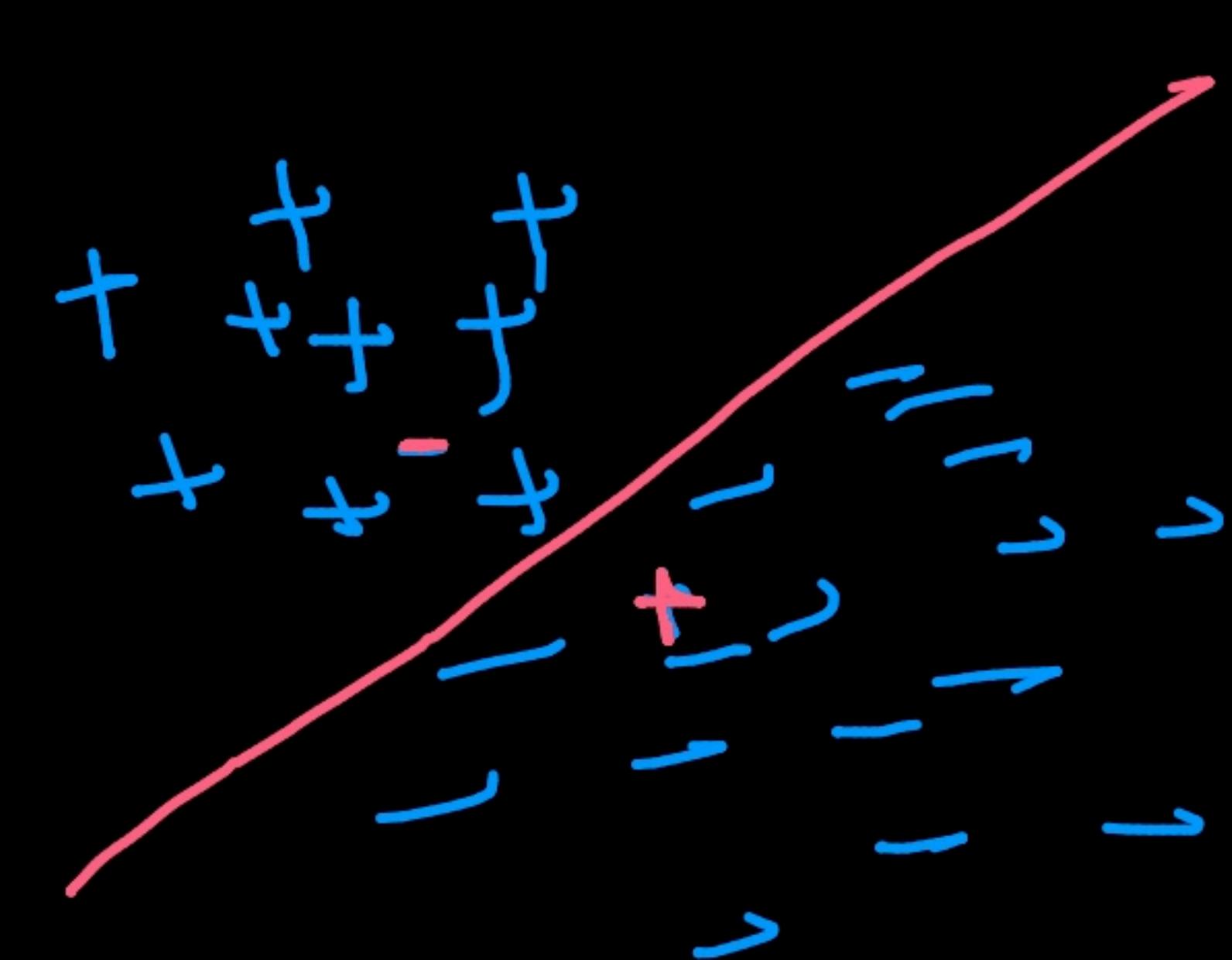
PCA







non-linearly Sep



+ Code + Text

\*\*\*\*\*  
2m ✓ 99936  
-6.907728433272202  
0.0010000368467572279  
\*\*\*\*\*  
{x} 99937  
-6.907738433540662  
0.0010000268460702846  
\*\*\*\*\*  
99938  
-6.907748433709118  
0.0010000168455833601  
\*\*\*\*\*  
99939  
-6.9077584337775715  
0.0010000068452964478  
\*\*\*\*\*  
-6.907768433746024  
0.0009999868453226368

