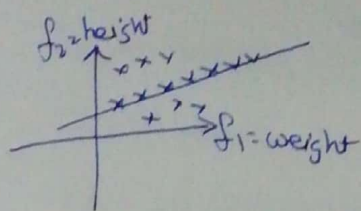


25.1 Geometric intuition of Linear Regression

Eg: Predict height given weight, gender, ethnicity, haircolor.

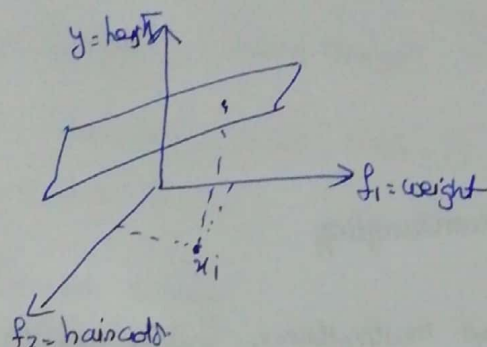
Lin Reg: find a line that fits the given data



$$y = w_1 * \text{weight} + w_0$$

$$y = mx + c$$

3 dimensional space

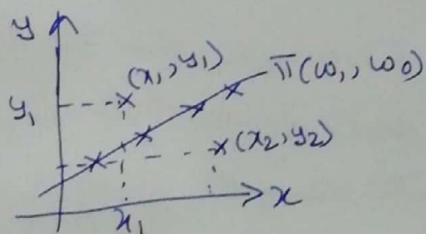


$$\text{height} = w_1 x_1 + w_2 x_2 + w_0$$

$$y_i = w_1 x_{i1} + w_2 x_{i2} + w_0$$

$$y = w^T x_i + w_0$$

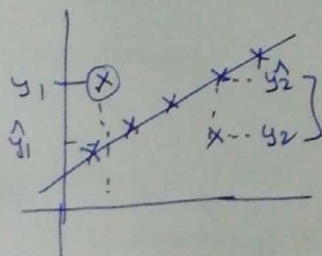
find a line/plane that best fits the datapoints



$$\text{Error} = y_1 - \hat{y}_1$$

$$\text{Error} = y_2 - \hat{y}_2$$

25.2 Mathematical Formulation



$$\text{Error}_1 = y_1 - \hat{y}_1 = +ve$$

$$\text{Error}_2 = y_2 - \hat{y}_2 = -ve$$

$$\text{Error}_3 = 0$$

we have to take a square

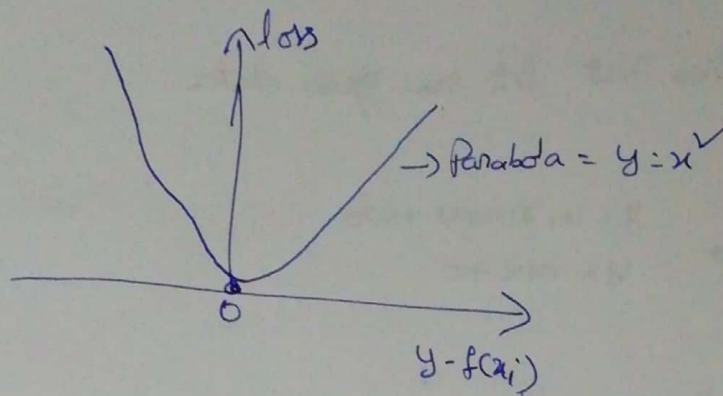
we have to find the optimal w^*, w_0^* =

$$\argmin_{w, w_0} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \hat{y}_i = f(x_i) = w^T x_i + w_0$$

$$(w^*, w_0^*) = \argmin_{w, w_0} \sum_{i=1}^n \{y_i - (w^T x_i + w_0)\}^2$$

Regularization

$$(w^*, w_0^*) = \underset{w, w_0}{\operatorname{argmin}} \sum_{i=1}^n \overbrace{\{y_i - (w^T x_i + w_0)\}^2}^{\text{square loss}} + \underbrace{\lambda \|w\|_2^2}_{L_2\text{-reg}}$$



25.3 Real world cases

Imbalanced data: upsampling & downsampling

Feature Importance & Interpretability: features are not multicollinear, feature weights can be used.

Feature Engineering
→ feature transformation:

Reg: L_1 reg, $\lambda_1 \uparrow$

Outliers: ~~log-reg~~ : sigmoid function → limiting the impact of outliers

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

~~sq loss~~ is highly impacted by outliers

RANSAC (Random Sampling)

- $D_{\text{train}} \rightarrow w^*, w_0^*$
- find pts very far away for π $(y_i - \hat{y}_i) \uparrow$
- remove these pts as outliers
- $D_{\text{train}} = D_{\text{train}} - \text{outlier}$
- repeat if needed

25.4 code sample for Linear Regression.

```
from sklearn.datasets import load_boston  
boston = load_boston()
```

```
// (506, 13)
```

```
boston.data DESCR
```

```
import pandas as pd
```

```
bos = pd.DataFrame(boston.data)
```

```
print(bos.head(1))
```

```
bos['Price'] = boston.target
```

```
X = bos.drop('PRICE', axis=1)
```

```
Y = bos['price']
```

```
import sklearn
```

```
// train & split use the code
```

```
from sklearn.linear_model import LinearRegression
```

```
lm = LinearRegression()
```

```
lm.fit(X_train, Y_train)
```

```
Y_pred = lm.predict(X_test)
```

```
plt.scatter(Y_test, Y_pred)
```

```
plt.xlabel("Prices: $ Y_i $")
```

```
plt.ylabel("Predicted Prices")
```

```
plt.show()
```

```
delta_y = Y_test - Y_pred
```

```
import seaborn as sns
```

```
import numpy as np
```

```
sns.set_style('white grid')
```

```
sns.kdeplot(np.array(delta_y), bw=0.5)
```

```
plt.show()
```

26.1 Differentiation

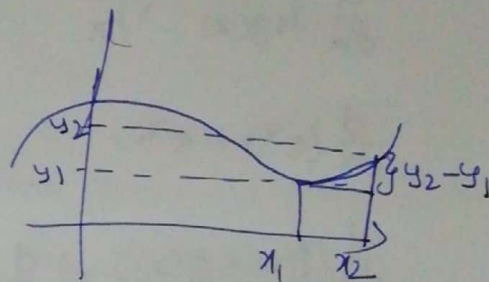
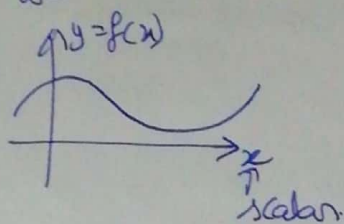
Solving optimization problems: differentiation.

ML \rightarrow differentiation \rightarrow Single value, Vector

\hookrightarrow maxima & minima.

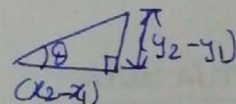
Single Variable diff:

$$y = f(x)$$



$$\left\{ \frac{dy}{dx} = \frac{df}{dx} = y' = f' \right\}$$

\uparrow differentiation of y w.r. to x .



$\frac{dy}{dx} \rightarrow$ Intuitively

\hookrightarrow Rate of change of y as x changes.

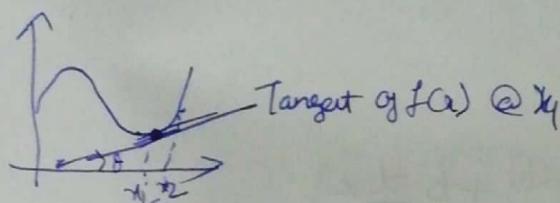
\hookrightarrow how much does y change as x changes.

$$\frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta y}{\Delta x} = \tan \theta$$

θ : Slope

$$\frac{dy}{dx} = \lim_{x \rightarrow 0} \frac{\Delta y}{\Delta x}$$

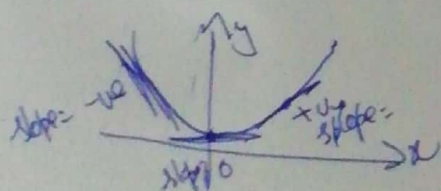
when x_2 comes closer to x_1 , then Δx will be almost closer to zero



Tangent is hypotenuse that we obtain as $\Delta x \rightarrow 0$

$$\text{Tang} = \frac{dy}{dx}$$

$\left[\frac{dy}{dx} \right]_{x_1}$: slope of the tangent to $f(x)$ @ $x = x_1$



basic

$$\frac{d}{dx}(x^n) = nx^{n-1}$$

$$\frac{d}{dx}(x^2) = 2x$$

$$\frac{d}{dx}(c) = 0$$

constant

$$\frac{d}{dx}(3x) = 3$$

$$\frac{d}{dx}(cx^n) = cnx^{n-1}$$

$$\frac{d}{dx} \log(x) = 1/x$$

$$\frac{d}{dx}(e^x) = e^x$$

$$\frac{d}{dx}(f(x) + g(x)) = \frac{d}{dx} f(x) + \frac{d}{dx} g(x)$$

chain rule

$$\frac{d}{dx} f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$f(g(x)) = (a-bx)^2$$

$$g(x) = a-bx$$

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

$$\frac{d}{dx} f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$\frac{dg}{dx} = \frac{d}{dx}(a-bx)$$

$$\begin{aligned} &= \frac{d}{dx}(a) - \left(x \frac{d}{dx}(b)\right) + b \frac{d}{dx}(x) \\ &= 0 \quad 0 \quad \downarrow \\ &= -b \quad 1 \end{aligned}$$

$$\frac{df}{dg}$$

$$g(x) = z$$

$$f(g(x)) = z^2$$

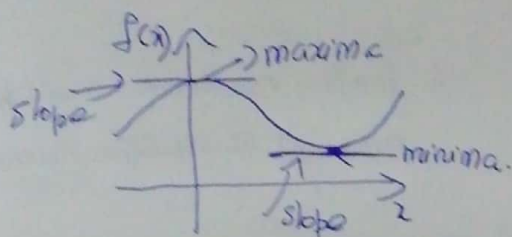
$$\frac{df}{dg} = \frac{dz^2}{dz} = 2z$$

$$\frac{d}{dx} f(g(x)) = -2b(a-bx)$$

26.3 Online differentiation tools

<https://www.derivative-calculator.net>

26.4 Maxima & minima



At minima and maxima slope will become zero

Slope is nothing but tan of angle b/w tangent & x axis.

$$f(x) = x^2 - 3x + 2 \quad (\text{find maxima \& minima})$$

$$\frac{df}{dx} = 0 \Rightarrow \text{slope} = 0$$

$$\frac{df}{dx} = 2x - 3 = 0 \Rightarrow \boxed{x = 3/2} \quad \text{At } x = 3/2 \text{ slope} = 0$$

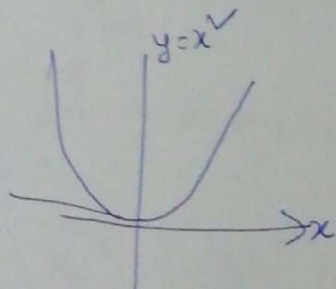
(Q) Is this a maxima & minima.

$$f(1.5) = (1.5)^2 - 3(1.5) + 2 = -0.25$$

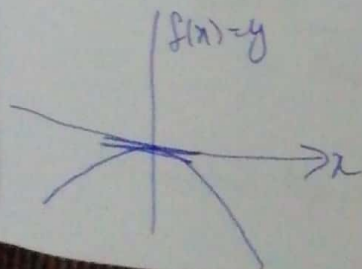
$$f(1) = 0 \quad \leftarrow \text{take a value close to } x.$$

$$f(1.5) < f(1)$$

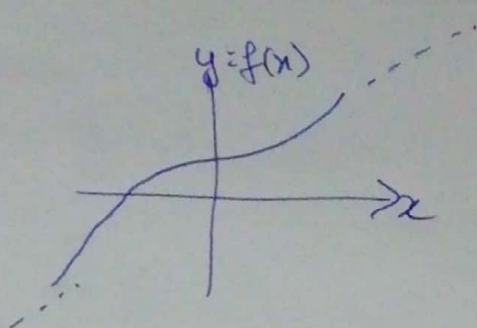
$f(1.5)$ cannot be maxima, it has to be minima



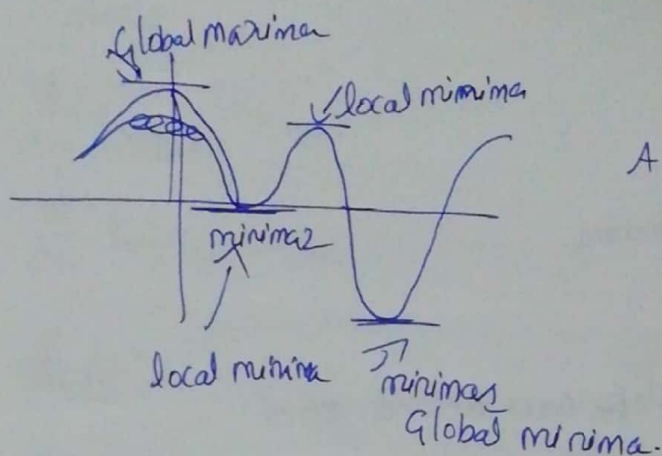
\rightarrow no maxima, it has minima at $x=0$



\rightarrow It has maxima but no minima



There is no maxima & no minima because as the dotted lines imply they can go to maxima & minima.



A function can have multiple minimal & multiple maxima.

minima of all the maxima is called global minima

Imagine

$$f(x) = \log(1 + \exp(ax))$$

find minima & maxima.

$$\frac{df}{dx} \frac{a \cdot \exp(ax)}{1 + \exp(ax)} = 0 \left\{ \begin{array}{l} \text{solving this is not trivial \& hard} \\ \text{To solve this we have computerised} \\ \text{program called } \boxed{\text{Gradient descent}} \\ \text{to solve maxima \& minima} \end{array} \right.$$

26.5 Vector Calculus : Grad.

$$y = a^T x$$

$$x = \langle x_1, x_2, \dots, x_d \rangle$$

$$a = \langle a_1, a_2, \dots, a_d \rangle \text{ Constant}$$

$$f(x) = y = \sum_{i=1}^d a_i x_i$$

$$\underbrace{\frac{df}{dx}}_{\text{vector}} = \underbrace{\nabla_x f}_{\text{grad or dot}}$$

$$\nabla_x f = \begin{bmatrix} \frac{df}{dx_1} \\ \frac{df}{dx_2} \\ \vdots \\ \frac{df}{dx_d} \end{bmatrix} \begin{matrix} \text{gradient} \\ \text{also is} \\ \text{vector} \end{matrix} \in \mathbb{R}^d$$

$$\frac{df}{dx_1} = \frac{\partial f}{\partial x_1} \rightarrow \text{Partial differentiation}$$

$$f(x) = y = a^T x = \sum_{i=1}^d a_i x_i$$

$$\nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} : \begin{matrix} \frac{\partial f}{\partial x_1} = a_1 \\ \frac{\partial f}{\partial x_d} = a_d \end{matrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{bmatrix}$$

$$\nabla_x (a^T x) = a$$

$$\underset{\text{logit loss}}{J(w)} = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda w^T w$$

x_i & y_i are constants. \leftarrow Data

$$\nabla_w J = \underbrace{-y_i x_i \frac{\exp(-y_i w^T x_i)}{1 + \exp(-y_i w^T x_i)}}_{\text{derivative}} + 2\lambda w = 0$$

Notice this Equation is not very easy

26.6 Gradient descent: geometric intuition

when dealing with maxima & minima

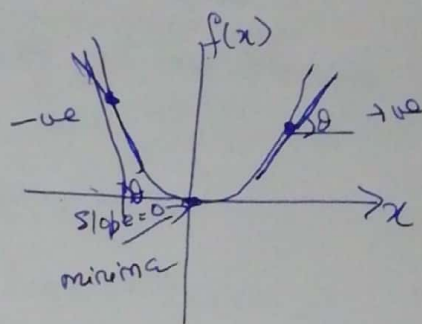
$$\frac{df}{dx} = 0 \quad ; \quad \nabla_x f = 0$$

→ Gradient Descent

↳ Iterative algorithm.

↳ pick a random number $x_0 \leftarrow$ first guess of x^*

↳ using gradient descent we will move to another point x_1 .
Eventually we will move closer to x^*



$$x^* = \underset{x}{\operatorname{argmin}} f(x)$$

$$\min f(x) = \max -f(x) \quad (8)$$

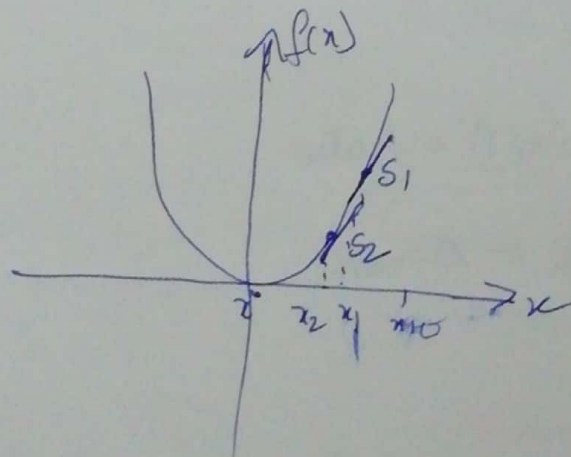
$$\max f(x) = \min -f(x)$$

minima

one side: slope = +ve

other side: slope = -ve.

→ slope changes its sign from +ve to -ve @ minima



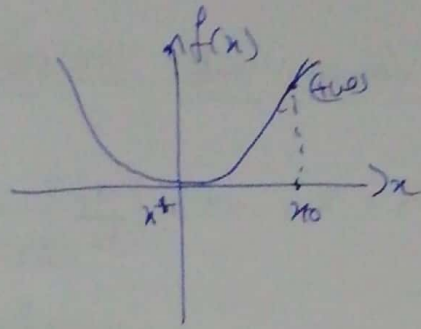
Approaching minima \leftarrow ^{Downside} slope decreases

Approach minima \rightarrow slope increases

Gradient Descent

① Pick an initial pt $= x_0$ @ random

② we have to pick x_1 such that x_1 is closer to x^* than x_0



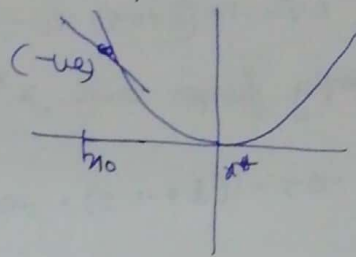
$$x_1 = x_0 - \gamma \left[\frac{df}{dx} \right]_{x_0}$$

Stepsize.

Let step-size $= 1$

$$x_1 = x_0 - 1 (+ve) \Rightarrow x_1 < x_0$$

~~②~~ $x_1 = x_0 - \gamma \left[\frac{df}{dx} \right]_{x_0}$
(=) low
 -ve



$$x_1 = x_0 - 1 (-ve)$$

$$\boxed{x_1 > x_0}$$

③ $x_2 = x_1 - \gamma \left[\frac{df}{dx} \right]_{x_1}$

At Any iteration

$$x_{i+1} = x_i - \gamma \left[\frac{df}{dx} \right]_{x_i}$$

$$x_0, x_1, x_2 \dots x_k$$

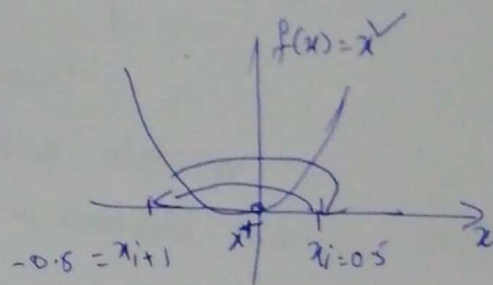
$$\boxed{x_k = x_{k-1} - \gamma \left[\frac{df}{dx} \right]_{x_{k-1}}} \leftarrow \text{update function.}$$

If $(x_{k+1} - x_k)$ is Very small

Then terminate @ $x^* = x_k$

26.7 Learning rate (or) step size

$$x_i = x_{i-1} - \gamma * \left[\frac{df}{dx} \right]_{x_{i-1}}$$



$$x_{i+1} = x_i - \gamma * \left[\frac{df}{dx} \right]_{x_i} \quad (\gamma = 1)$$

$$x_{i+1} = 0.5 - 1 * (2 * 0.5) = -0.5$$

we simply jump over x^*

$$x_{i+2} = -0.5 - 1 * (2 * -0.5) = -0.5 - 1 * (-1) = 0.5$$

Remedy for Oscillation:

→ Change γ with Each iteration

→ Reduce γ for Each iteration

s.t. iterations \uparrow and $\gamma \downarrow$

26.8 Gradient descent for linear regression

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (y_i - w^T x_i)^2$$

\uparrow no-bias
 \uparrow allowing w_0 (intercept)

$$L = \sum_{i=1}^n (y_i - w^T x_i)^2$$

$\because x_i \& y_i$ are constants
 \uparrow Vector \uparrow scalar

$$\nabla_w L = \sum_{i=1}^n \{ 2(y_i - w^T x_i)(-x_i) \}$$

Pick a random vector $w_0 = \langle \dots \rangle$

$$w_1 = w_0 - \gamma \sum_{i=1}^n (-2x_i)(y_i - w_0^T x_i)$$

$$w_2 = w_1 - \gamma \sum_{i=1}^n (-2x_i)(y_i - w_1^T x_i)$$

We will calculate $w_0, w_1, w_2, \dots, w_k, w_{k+1}$

when $w_{k+1} - w_k$ is very small then we will declare $w^* = w_k$

Problem

In n is very large $n = 1 \text{ million}$, computing summation is

Very Expensive. Gradient descent will become very slow

Solution

Stochastic Gradient descent is the solution.

26.9 SGD algorithm

Linear regression:

$$w_{j+1} = w_j - \gamma \sum_{i=1}^n (-2x_i)(y_i - w_j^T x_i)$$

SGD (Stochastic Gradient Descent)

$$w_{j+1} = w_j - \eta \sum_{i=1}^K (-x_i) (y_i - w_j^T x_i)$$

$$1 \leq i \leq n$$

→ Pick a random set of K -pts

$$w_{GD}^* = w_{SGD}^*$$

GD: ($n=1$ million) 100 iterations to converge x^*

SGD: $n=1000$, in >100 iterations x^*

SGD: $n=10$; in >1000 iteration to converge at x^*

K = # of points pick @ iteration.

Every iterations points in K is different

K = batch size in SGD

batch of random pts

SGD: The most imp-optimization algo in ML

26.16 Constrained Optimization \rightarrow PCA

Till now we have seen optimization problems like ~~max~~
 $\max_x f(x)$ or $\min_x f(x)$

If we recall PCA \rightarrow objective.

$$\max_u \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

$$\text{s.t. } u^T u = 1$$

\downarrow
Constraint

$$\approx \max_n f(x)$$

$$\text{s.t. } g(x) = c$$

general Constraint optimization.

$$\max_x f(x) \text{ objective function.}$$

$$\text{s.t. } g(x) = c \rightarrow \text{Equality Constraint}$$

$$h(x) \geq d \rightarrow \text{inequality Constraint}$$

$$k(x) \leq e \rightarrow -k(x) \geq -e$$

can be converted by multiplying both side with negative sign.

Lagrangian multipliers

$$\mathcal{L} = f(x) - \lambda \{g(x) - c\} - \mu \{d - h(x)\}$$

λ & μ ~~are~~ are called Lagrangian multipliers.

$$\lambda \geq 0, \mu \geq 0$$

$$\frac{\partial \mathcal{L}}{\partial x} = 0; \quad \frac{\partial \mathcal{L}}{\partial \lambda} = 0; \quad \frac{\partial \mathcal{L}}{\partial \mu} = 0$$

$$\downarrow$$

 \tilde{x}

$$\downarrow$$

 $\tilde{\lambda}$

$$\downarrow$$

 $\tilde{\mu}$

$$x^* = \max_x f(x)$$

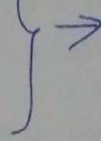
Here

$$\tilde{x} = x^*$$

PCA:

$$\max_u \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

$$s.t. \quad u^T u = 1$$



$$\max_u \quad u^T S u$$

$$s.t. \quad u^T u = 1$$

$$S = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$



Covariance matrix of X

$$\mathcal{L}(u, \lambda) = u^T S u - \lambda (u^T u - 1)$$

$$\frac{\partial \mathcal{L}}{\partial u} = 0 \Rightarrow \frac{\partial}{\partial u} (u^T S u - \lambda u^T u - \lambda) = 0$$

$$= 2Su - \lambda u = 0$$

$$\Rightarrow Su = \lambda u$$

Covariance
matrix

↑
Vecto.

26.11 Logistic Regression formulation revisited

$$\omega^* = \underset{\omega}{\operatorname{argmin}} (\text{logistic-loss}) + \text{regularization} (\lambda \omega^T \omega)$$

Lagrangian

$$\begin{aligned} \mathcal{L} &= \text{logistic-loss} - \lambda (1 - \omega^T \omega) \\ &= \text{logistic-loss} - \lambda + \lambda \omega^T \omega \end{aligned}$$

regularization can be thought as imposing an Eq. constraint.

26.12 Why L1 regularization creates sparsity?

Logistic regression + L1 reg creates sparsity in ω as compared to L2 regularization.

Sparsity means: $\omega \in \langle \omega_0, \omega_1, \dots, \omega_d \rangle$

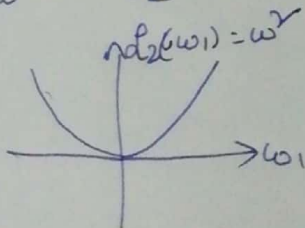
most of ω_i 's are "0"

\mathcal{L}_2

$$\min_{\omega} \text{loss} + \lambda \|\omega\|_2^2$$

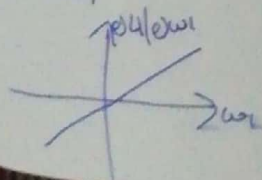
$$\min_{\omega_1, \omega_2, \dots, \omega_d} (\omega_1^2 + \omega_2^2 + \dots + \omega_d^2)$$

$$\min_{\omega} \omega^2 = \mathcal{L}_2(\omega)$$



$$\mathcal{L}_2(\omega_1) = \omega_1^2$$

$$\frac{\partial \mathcal{L}_2}{\partial \omega_1} = 2\omega_1$$

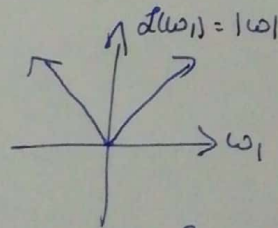


\mathcal{L}_1

$$\min_{\omega} \text{loss} + \lambda \|\omega\|_1$$

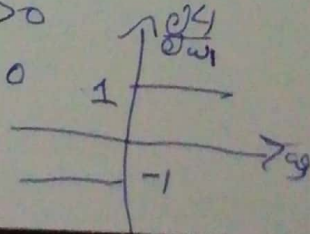
$$\min_{\omega_1, \omega_2, \dots, \omega_d} (|\omega_1| + |\omega_2| + \dots + |\omega_d|)$$

$$\min_{\omega} |\omega| = \mathcal{L}_1(\omega)$$



$$\mathcal{L}_1(\omega_1) = |\omega_1| = \begin{cases} \omega_1 & \text{if } \omega_1 > 0 \\ -\omega_1 & \text{if } \omega_1 \leq 0 \end{cases}$$

$$\frac{\partial \mathcal{L}_1}{\partial \omega_1} = \begin{cases} +1 & \text{if } \omega_1 > 0 \\ -1 & \text{if } \omega_1 < 0 \end{cases}$$



\mathcal{L}_2

$$(w_1)_{j+1} = (w_1)_j - \eta \left[\frac{\partial \mathcal{L}_2}{\partial w_1} \right]_{w_{1j}}$$

Let $(w_{1j}) = +ve$

$$(w_1)_{j+1} = (w_1)_j - \gamma(2w_{1j})$$

 \mathcal{L}_1

$$(w_1)_{j+1} = (w_1)_j - \eta \left(\frac{\partial \mathcal{L}_1}{\partial w_1} \right)_{w_{1j}}$$

Let $(w_{1j}) = +ve$

$$(w_1)_{j+1} = (w_1)_j - \gamma(1)$$

26.13 Implement SGD for Linear Regression