

AI vs ML

Data is knowledge. AI is the representation & manipulation of that knowledge.

Deep Learning is a branch of Neural Network.

"How much does a rhino weight?"

Descriptive Statistic : mean / median / mode.

Is the data we have **valid enough** to answer?

Did not miss a rhino.

PDF. x has pdf f . $\int_a^b f(x) dx$: $P(a \leq x \leq b)$

cdf

$$F(x) = \int_{-\infty}^x f(e) de = P(X \leq x)$$

Average weight of rhino.

True answer.

$$\boxed{\mu} = E[X] = \int_{-\infty}^{\infty} x f(x) dx$$

Procedure :

$$\bar{x}_n = \frac{1}{n} (x_1 + \dots + x_n)$$

How close is \bar{x}_n to μ ?

$$\begin{aligned}
 E[\bar{X}_n] &= E\left[\frac{1}{n}(X_1 + \dots + X_n)\right] \\
 &= \frac{1}{n}E[X_1 + \dots + X_n] \\
 &= \frac{1}{n}(E[X_1] + \dots + E[X_n]) \\
 &= \frac{1}{n}(\mu + \dots + \mu) = \underline{\mu} \quad \text{valid way}
 \end{aligned}$$

Variance.

$$\text{Var}(x) = E[(x - \mu)^2] \Rightarrow \text{expect } \text{Var}(x) = \sigma^2 \text{ to be zero.}$$

$$E[x] = \mu$$

Variance (\bar{X}_n)

$$\begin{aligned}
 &= \text{Var}\left(\frac{1}{n}(X_1 + \dots + X_n)\right) \\
 &= \frac{1}{n^2} \text{Var}(X_1 + \dots + X_n) \\
 &= \frac{1}{n^2} [\text{Var}(X_1) + \dots + \text{Var}(X_n)] \\
 &= \frac{1}{n^2} [\sigma^2 + \dots + \sigma^2] \\
 &= \frac{\sigma^2}{n} \quad \Rightarrow \text{Var}(\bar{X}_n) = E[(\bar{X}_n - \mu)^2] = \frac{\sigma^2}{n}
 \end{aligned}$$

$n \uparrow \text{Var}(\bar{X}_n) \downarrow$ more accurate.

"How confident can we be in our estimate?"

$$\begin{aligned}
 P(|\bar{X}_n - \mu| \geq \varepsilon) \\
 \leq \frac{\text{Var}(\bar{X}_n)}{\varepsilon^2} = \frac{\sigma^2}{\varepsilon^2} = O\left(\frac{1}{n}\right)
 \end{aligned}$$

Chebyshev.

$$P(|x - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2}$$

The more we measure, the more likely

it falls in this range.



$$X_i \in \{0, 1\} \quad E[X_i] = P(X_i=0) \cdot 0 + P(X_i=1) \cdot 1 \\ = P(X_i=1) = \mu$$

Hoeffding Inequality.

$$P(|\bar{X}_n - \mu| \geq \varepsilon) \leq 2e^{-2\varepsilon^2 n}$$

(only hold for binary measuring)

$$P(|\bar{X}_n - \mu| \geq \varepsilon) \leq \delta$$

$$2e^{-2\varepsilon^2 n} \leq \delta$$

$$e^{-2\varepsilon^2 n} \leq \frac{1}{2}\delta$$

$$-2\varepsilon^2 n \leq \ln(\frac{1}{2}\delta)$$

$$n \geq -\frac{1}{2\varepsilon^2} \ln(\frac{1}{2}\delta)$$

$$n \geq \frac{1}{\varepsilon^2} \ln(\frac{2}{\delta})$$

$$n \geq O\left(\frac{1}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right)\right)$$

$\varepsilon \downarrow \delta \downarrow$ and $n \uparrow$ (more sample)

For binary data, we need less sample.

Maximum Likelihood Estimation.

What parameter make your data as likely as possible?

Data: x_1, x_2, \dots, x_n

$$\begin{aligned}lik_{\mu\sigma^2}(\text{data}) &= f(x_1) \cdot f(x_2) \cdots f(x_n) \\&= \prod_{i=1}^N \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \\&= \frac{1}{\sigma^N} (2\pi)^{\frac{N}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2\right) \\&= \ln\left(\frac{1}{\sigma^N} (2\pi)^{\frac{N}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2\right)\right) \\&= -N \ln(\sigma) - \frac{N}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2\end{aligned}$$

$$L = \ln lik_{\mu\sigma^2}(\text{data}) = -N \ln(\sigma) - \frac{N}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

$$\begin{aligned}\frac{\partial L}{\partial \mu} &= 0 - 0 - \frac{1}{2\sigma^2} \sum_{i=1}^N -2(x_i - \mu) = 0 \\&= \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) = 0\end{aligned}$$

find a extremum
of a function

take derivative and
set to zero.

$$\frac{\partial L}{\partial \sigma} = -\frac{N}{\sigma} - \frac{(-2)}{2\sigma^3} \sum_{i=1}^N (x_i - \mu)^2 = 0$$

$$\sum_{i=1}^N (x_i - \mu) = 0 \quad -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^N (x_i - \hat{\mu})^2 = 0$$

$$\left(\sum_{i=1}^N x_i\right) - \left(\sum_{i=1}^N \mu\right) = 0 \quad \frac{1}{\sigma^3} \sum_{i=1}^N (x_i - \hat{\mu})^2 = \frac{N}{\sigma^2}$$

$$\sum_{i=1}^N x_i - N\mu = 0$$

$$\frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2 = \sigma^2$$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \leftarrow \text{estimate of true value}$$

$$X \sim \text{Unif}[0, L]$$

$$\text{Uniform Distribution. } E[x] = \int_R x f(x) dx$$

pdf

$$X_1, X_2, \dots, X_N \sim \text{Unif}[0, L]$$

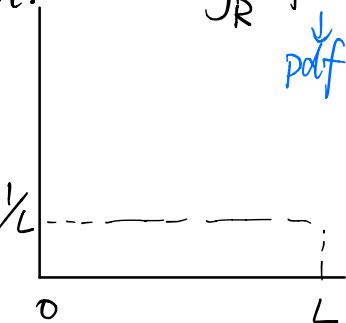
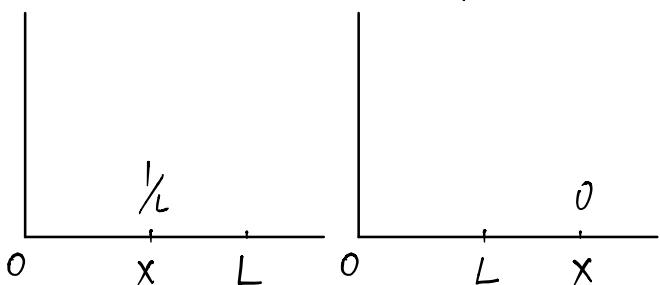
From data, estimate L .

M.O.M

$$E[x] = \frac{1}{2}L$$

$$\hat{L} = 2\bar{x}_n$$

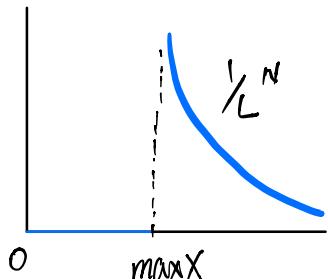
$$\text{likelihood}_L(x) = \begin{cases} \frac{1}{L} & L \geq x \\ 0 & \text{else.} \end{cases}$$



$$\begin{aligned} \int_0^L x \left(\frac{1}{L}\right) dx &\rightarrow \text{pdf of uniform Distribution} \\ &= \frac{1}{L} \cdot \frac{1}{2}x^2 \Big|_0^L \\ &= \frac{1}{L} \cdot \frac{1}{2} \cdot L^2 = \frac{1}{2}L \end{aligned}$$

$$\text{lik}_L(\text{data}) = \begin{cases} \left(\frac{1}{L}\right)^N & L \geq \max X \\ 0 & \text{else} \end{cases}$$

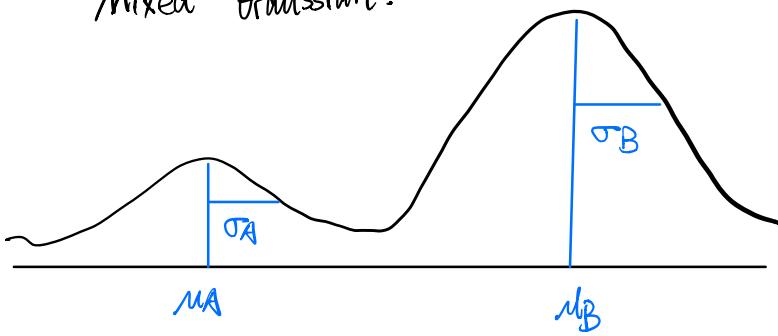
Find L to maximise lik_L .



$$\hat{L}_{MLE} = \max \bar{x}_i$$

$$\begin{aligned} P(|\hat{L}_{MLE} - L| \geq \varepsilon) \\ P(|\hat{L}_{MLE} - L| \geq \varepsilon) \end{aligned}$$

Mixed Gaussian.



Typicality of Data

(last 10 days. lost 7 times.)

$H_0:$
 hypothesis: the coin is fair
 $H_1:$
 hypothesis: the coin is not fair
 (The coin is biased toward tails)

$$P(7 \text{ tails in 10 flips} | H_0)$$

$$P(7 \text{ or more tails in 10 flips} | H_0)$$

$$= \sum_{k=7}^{10} \binom{10}{k} \left(\frac{1}{2}\right)^k \left(\frac{1}{2}\right)^{10-k} = 0.17$$

$$\leq 0.05$$

$$P(70 \text{ tails in 100 flips} | H_0)$$

$$= \sum_{k=70}^{100} \binom{100}{k} \left(\frac{1}{2}\right)^k \left(\frac{1}{2}\right)^{100-k} = 0.000039$$

extremely unlikely

Pill A / 50%

Pill B / 70%

Bayesian Analysis

$H_0:$ the pills are pill A
 $H_1:$ the pills are pill B

How much should you believe
the hypothesis given
the data.

$P(\text{Data} | H_0) \Rightarrow P(H_0 | \text{Data})$
 $P(\text{Data} | H_1) = \frac{P(H_0) P(\text{Data} | H_0)}{P(\text{Data})}$
 $= \frac{P(H_0) P(\text{Data} | H_0)}{P(H_0) P(\text{Data} | H_0) + P(H_1) P(\text{Data} | H_1)}$
 $P(H_1 | \text{Data}) = \frac{P(H_1) P(\text{Data} | H_1)}{P(H_0) P(\text{Data} | H_0) + P(H_1) P(\text{Data} | H_1)}$

if $P(H_0) = P(H_1) = \frac{1}{2}$

$P(H_0 | \text{Data}) = \frac{P_0}{P_0 + P_1}$
 $P(H_1 | \text{Data}) = \frac{P_1}{P_0 + P_1}$

Hypothesis. X, Y are independent.

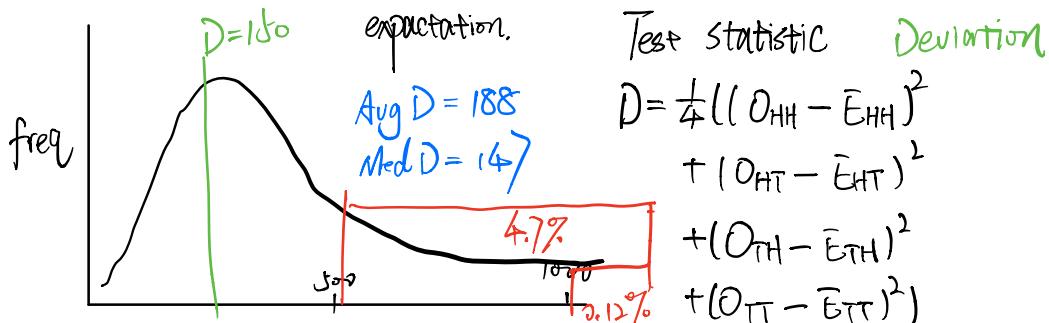
flip coin X, coin Y 1000 times together.

$$O_{HH} = 236 (28\%)$$

$$O_{HT} = 265 (24\%) \quad P(X=H) = \frac{1}{2}$$

$$O_{TH} = 240 (25\%) \quad P(Y=H) = \frac{1}{2}$$

$$O_{TT} = 259 (28\%)$$



$$\bar{E}_{HH} = \overset{\rightarrow}{P(X=H)} \cdot \overset{\leftarrow}{P(Y=H)} \cdot N = \frac{N}{4}$$

$$P = \frac{O_{HT} + O_{HH}}{O_{HH} + O_{HT} + O_{TH} + O_{TT}}$$

$$Q = \frac{O_{HT} + O_{TH}}{O_{HH} + O_{HT} + O_{TH} + O_{TT}}$$

$$\bar{E}_{HH} = PQ, N \quad T = \sum_x \sum_y \frac{(O_{xy} - \bar{E}_{xy})^2}{\bar{E}_{xy}}$$

If x, y are independent, then $T \sim \chi^2_1, P(T \geq T_0)$

$$= \int_{T_0}^{\infty} \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{N}} e^{-\frac{t^2}{2}} dt$$

$$P(T \geq 3.84) = 5\%$$

X has r values. Y has c values

$$T \sim \chi^2_{(r-1)(c-1)}$$

Data
Input $x = (x_1, x_2, \dots, x_d)$
Output $y = 0/1$

$\{(x_i, y_i)\}_{i=1}^m$ i.i.d from some distribution.
Given $x=x$, find $P(Y=y | X=x)$.

Given $\{(x_i, y_i)\}_{i=1}^m$. Assume $y = f(x)$ for some f .

How can we find f to model our data?

Hypothesis space $f: x \rightarrow y$. Set of boolean function

find hypothesis f

x_1	x_2	y
1	1	0
1	0	1
0	1	1
0	0	0

Truth table.

$$\text{err}(f) = \overline{P}(f(x) \neq y)$$

Goal: find f to minimize $\text{err}(f)$.

↓ we don't know, so we can not compute $\text{err}(f)$.

This goal doesn't stand.

if $f(x_i) \neq y_i$, = 1

else = 0.

$$\text{err}_{\text{train}}(f) = \frac{1}{m} \sum_{i=1}^m \underbrace{\mathbb{1}\{f(x_i) \neq y_i\}}$$

= avg missclassification on the data

$$\text{Solve: } \min_{f \in H} \text{err}_{\text{train}}(f) \Rightarrow f^*$$

f is overfit if $\text{err}_{\text{train}}(f)$ small, $\text{err}(f)$ large.

if $f^* \rightarrow$ true f . $\hat{f} \rightarrow$ estimate f .

$$\hat{f}(x) = \begin{cases} y_i & \text{if } x=x_i \\ 1-f^*(x) & \text{else} \end{cases}$$

estimate \hat{f} is zero.
but f^* could be quite large.

If H is finite, if we want $|\text{err}(f) - \text{err}_{\text{train}}(f)| < \varepsilon$.

$$m \geq O\left(\frac{1}{\varepsilon^2} \ln |H|\right)$$

the more hypothesis, the more data we need.
number of functions in H .

H : boolean functions on D variables.

Truth table

	x_1	x_2	x_3	\dots	x_D		fix	$ H = 2^{12^D}$
2^D	\uparrow	0	$;$	$;$	$;$			
	$;$	$;$	$;$	$;$	$;$			

How many data points do we need to prevent overfitting?

$$\frac{1}{\varepsilon^2} \ln |H| = \frac{1}{\varepsilon^2} \ln(2^{12^D}) = O(2^D)$$

$$= 2^D \frac{1}{\varepsilon^2} \ln(2)$$

Regularization: Restricting H .

- Occam's Razor.
- If H is smaller, we need less data to prevent overfitting.
- Simpler models have less room for spurious relationships.
假的

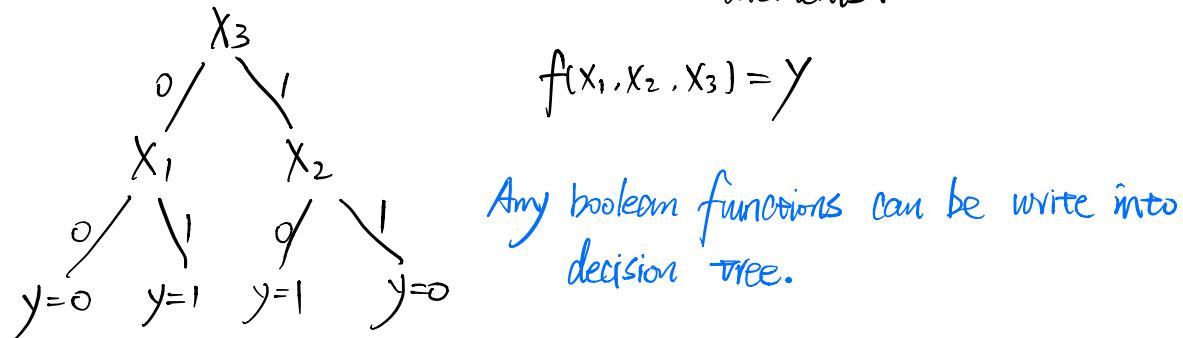
Problem of machine learning.

$$\min_{f \in H} \text{err}_{\text{train}}(f) \rightarrow \text{algorithms are trying to solve this problem.}$$

To estimate $\text{err}_{\text{train}}(f) \approx \text{err}(f)$

- read a lot of data
- need few hypothesis.

Decision Tree. Graphical models of restricted if/else statements.



- 1) How to fit decision trees to data?
- 2) How can we regularize?
 - Learn towards shorter trees.
- 3) Given a tree, how can we improve it?

Algorithm.

- 1) Pick a decision variable X_i .
- 2) Partition data : $x_i = 0/1$
- 3) Recursively split on each part.
- 4) Continue until all data in one branch has the same y value.

ID3.

X_1	X_2	X_3	Y
1	1	0	1
1	0	1	1
1	0	0	0
0	0	0	?

We don't know.
No data about X_1 .

entropy

$$\underline{H}(Y) = \sum_y -P(Y=y) \log_2 P(Y=y)$$

$$y = C. H(Y) = 0.$$

$$y = \{1, 2, 3, \dots, 1024\}. H(Y) = \sum_{i=1}^{1024} -\left(\frac{1}{1024}\right) \log_2 1024 \\ = 10 \times \frac{1}{1024} \times 1024 = 10$$

$$H(Y|X=x) = -\sum_y P(Y=y|X=x) \log_2 P(Y=y|X=x)$$

Suppose x, y are independent. $P(Y=y|X=x) = P(Y=y)$ X tell us nothing.

Aug information $\rightarrow H(y|x) = \sum_x P(x=x) H(Y|X=x)$
 x dependency over y

$$\underline{IG}(x) = H(y) - H(Y|x)$$

information game. if $x=Y$. $\underline{IG}(x) = H(Y) - H(Y|y) = H(Y)$

$$H(y|y) = \sum_y P(y=y) H(Y|y=y) = 0$$

x, y are independent. $H(Y|x) = H(Y)$

$$\underline{IG}(x) = H(Y) - H(Y|x) = 0.$$

2D3 - Iterative Dichotomiser 3

- 1) For each X_i , estimate $\underline{IG}(x_i)$.
- 2) Find X_i with maximum \underline{IG} and split on it.
- 3) Recurse on each part.
- 4) Stop when the remaining data having same Y .

Estimation of ID3.

$$IG(x_3) = H(Y) - H(Y|x_3)$$

Data Set

$$H(Y) = -P(Y=0) \log_2(P(Y=0))$$

$$- P(Y=1) \log_2(P(Y=1))$$

$$P(Y=0) = \frac{\#\{y=0\}}{\text{total \# samples.}}$$

\Downarrow estimation.

$$\hat{IG}(x_3) = \hat{H}(Y) - H(Y|x_3).$$

$$H(Y|x_3) = P(x_3=0)H(Y|x_3=0)$$

$$+ P(x_3=1)H(Y|x_3=1)$$

$$H(y|x_3=0) = -P(Y=0|x_3=0) \log_2 P(Y=0|x_3=0)$$

$$- P(Y=1|x_3=0) \log_2 P(Y=1|x_3=0)$$

$$P(Y=0|x_3=0) = \frac{P(Y=0, x_3=0)}{P(x_3=0)} = \frac{\#\{y=0, x_3=0\}/N}{\#\{x_3=0\}/N}$$

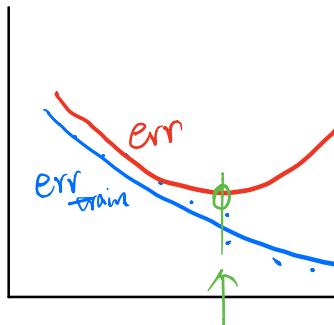
$$\hat{IG}(x_3) = \hat{H}(Y) - \hat{H}(Y|x_3)$$

How to prevent overfitting.

- Stop when we hit a certain depth.
- Assign by majority
- Assess probability.
- Stop when you hit certain sample size.

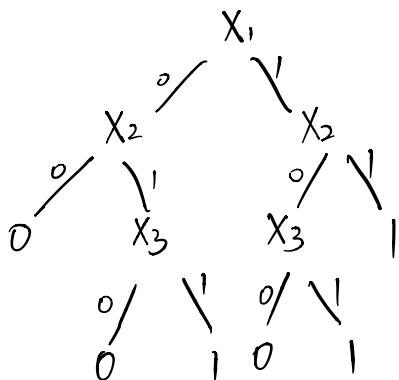
$\text{err}_{\text{train}}(f)$
 $\text{err}(f)$

Data: 80/20



terminate training

$$f(x_1, x_2, x_3) = \text{majority}(x_1, x_2, x_3)$$



$f = 0 \cdot \text{output} + 1$.

Perceptron

$$f(x) = \text{sgn}(w_1 x_1 + w_2 x_2 + \dots + w_D x_D + b)$$

$$f(x_1, x_2, x_3) = \text{majority}(x_1, x_2, x_3). \quad x_1, x_2, x_3 \in \{+1, -1\}$$

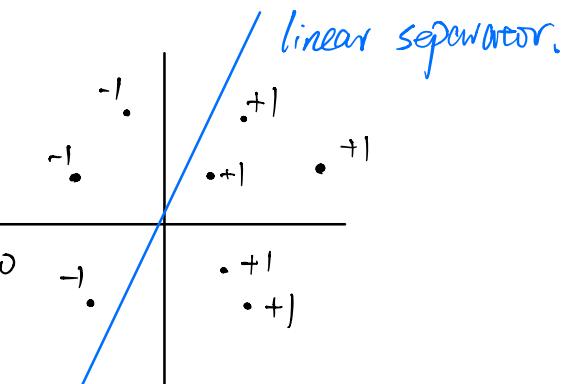
1-D data . $x \in \mathbb{R}$. $y \in \{+1, -1\}$

$$f(x) = \begin{cases} +1 & x \geq \frac{-b}{w} \\ -1 & \text{else.} \end{cases}$$

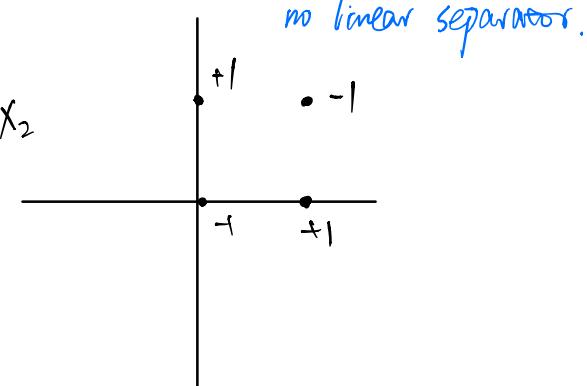
$$X = (x_1, x_2) \in \mathbb{R}^2$$

$$y \in \{-1, +1\}$$

$$f(x_1, x_2) = \begin{cases} +1 & w_1 x_1 + w_2 x_2 + b \geq 0 \\ -1 & \text{else} \end{cases}$$



$$f(x_1, x_2) = x_1 \text{ xor } x_2$$



$$f(\underline{x}) = \begin{cases} +1 & \underline{w} \cdot \underline{x} + b \geq 0 \\ -1 & \text{else} \end{cases}$$

Perceptron.

- w ^{Vector} = 0, $b = 0$ if no linear separator. gonna run forever.
- Find a misclassified point (\underline{x}, y)
- $\underline{w} \leftarrow \underline{w} + y \underline{x}$ $f_{\underline{w}}(\underline{x}) \neq y$
 $b \leftarrow b + y$
- Repeat

current model w, b

$$\underline{w}' = \underline{w} + y \underline{x}$$

$$b' = b + y$$

$$f(\underline{x}) = +1, y = -1$$

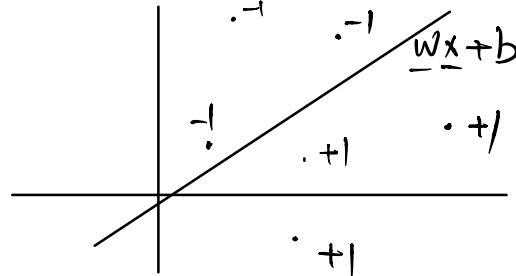
$$f(\underline{x}) = -1, y = +1$$

$$(\underline{w} \cdot \underline{x} + b) y < 0$$

$$\begin{aligned}
 & \underline{w}' \cdot \underline{x} + b' && \text{misclassified.} \\
 & = (\underline{w} + y \underline{x}) \cdot \underline{x} + (b + y) && (\underline{w} \cdot \underline{x} + b) \cdot y < 0 \\
 & = \underline{w} \cdot \underline{x} + b + y (||\underline{x}||^2 + 1) && \text{correctly classified.} \\
 & (\underline{w}' \cdot \underline{x} + b')y = (\underline{w} \cdot \underline{x} + b)y + y^2 (||\underline{x}||^2 + 1) && (\underline{w} \cdot \underline{x} + b) \cdot y > 0
 \end{aligned}$$

$\gamma(\underline{w}, b)$ = margin of the separator $\underline{w} \cdot \underline{x} + b$
 = the distance of the closest data point to $\underline{w} \cdot \underline{x} + b$

$$= \min_i \frac{|\underline{w} \cdot \underline{x}_i + b|}{\|\underline{w}\|}$$



If a linear separator exists ($\underline{w} \cdot \underline{x} + b = 0$).

the PLA terminates in at most $T \leq \left(\frac{2}{\gamma(\underline{w}, b)}\right)^2$ steps.

Run the algo certain times. if it outputs a linear separator,
 Great. If not, then either there's no linear separator for
 this data or the data is not reasonably separable.
 (too close. separating this is
 basically noise and not important).

Why perceptron?

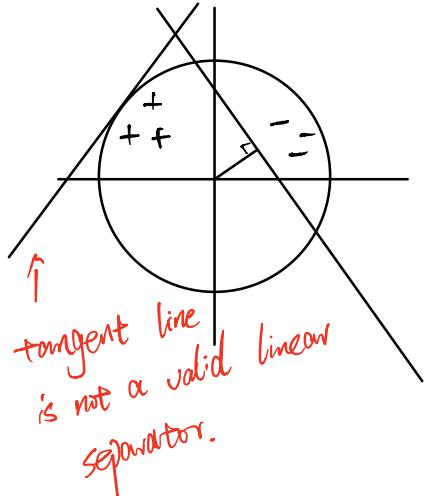
Vs. decision tree. different model but giving same output
 is closer geometrically.

PLA.

- $\underline{w} \cdot b = 0$
- find misclassified data \underline{x}, y
- update $\underline{w} \leftarrow \underline{w} + y\underline{x}$
 $b \leftarrow b + y$
- Repeat

$$f(\underline{w}, b) = \min_i \frac{|\underline{w} \cdot \underline{x}_i + b|}{\|\underline{w}\|_2}$$

Assume $\|\underline{x}_i\| \leq 1$. if $\underline{w} \cdot \underline{x} + b = 0$ cuts into the unit circle,
that, $b^2 / \|\underline{w}\|^2 \leq 1$. $b^2 \leq \|\underline{w}\|^2$



H : sets of all possible perceptions on a data set (from the PLA)
 $|H| \leq m^T \leq m^{\frac{4}{(\gamma^*)^2}}$ (size m).

Perceptions.

$$\text{Given } \{(x_i, y_i)\}_{i=1}^m \quad T \leq \left(\frac{2}{\gamma^*}\right)^2$$

$$W = \underline{x}_1 y_1 + \underline{x}_2 y_2 + \dots + \underline{x}_T y_T$$

$\xrightarrow{\text{iterations}}$

$$|H| \leq m^T \leq m^{\left(\frac{2}{\gamma^*}\right)^2}$$

γ^* maximum separation between two data sets.

$$|\text{err}_{\text{train}}(f) - \text{err}(f)| < \varepsilon \rightarrow m \geq \frac{1}{\varepsilon^2} \log |H|$$

$$m \geq \frac{1}{\varepsilon^2} \left(\frac{2}{\gamma^*}\right)^2 \log(m) \rightarrow m \geq \frac{4}{\varepsilon^2 (\gamma^*)^2} \log^2 \left(\frac{4}{\varepsilon^2 (\gamma^*)^2}\right)$$

Size of data

$$\gamma(\underline{w}, b) = \min_i \frac{|\underline{w}x_i + b|}{\|\underline{w}\|_2} \quad \gamma^* = \max_{\substack{\underline{w}, b \\ \text{valid}}} \gamma(\underline{w}, b)$$

\rightarrow valid linear separator to separate two data sets

Constructed
Satisfactory

Find \underline{w}, b s.t. $\nexists \text{sgn}(\underline{w}x_i + b) \neq y_i$

↓

Find \underline{w}, b s.t. $\forall y_i (\underline{w}x_i + b) > 0$.

$$w_1(x_1^i y_i) + w_2(x_2^i y_i) + \dots + w_b(x_b^i y_i) + b y > 0$$

Linear Programming

$$\max_{x_1, \dots, x_n} C_1 x_1 + \dots + C_n x_n . \text{ s.t. } Ax \geq b$$

To find a linear separator solve trying to map this problem into linear programming problem.

$$\begin{aligned} \max_{\underline{w}, b} \quad & b \\ \text{vector} \quad & \downarrow \\ y_i(\underline{w} \underline{x}_i + b) > \underline{\delta} & \text{strictly positive} \\ \downarrow \text{scale.} & \\ y_i(\underline{w} \underline{x}_i + b) \geq 1 & \end{aligned}$$

There must be a δ s.t. $\underline{\delta} > 0$ and $\underline{\delta} \geq \delta$.

$$m \geq \frac{1}{\varepsilon^2} d \quad \text{to prevent overfitting.}$$

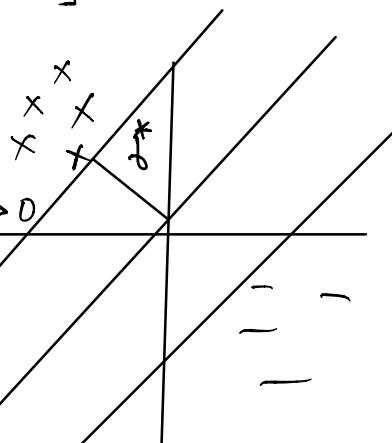
$$m \geq \min \left[\frac{4}{\varepsilon^2 (\gamma^*)^2} \log^2 \left(\frac{4}{\varepsilon^2 (\gamma^*)^2} \right), \frac{1}{\varepsilon^2} d \right]$$

Support Vector Machines

$$\max_{\underline{w}, b} \gamma(\underline{w}, b) . \text{ s.t. } \forall i : y_i(\underline{w} \underline{x}_i + b) \geq 0$$

$$\max_{\underline{w}, b} \frac{1}{\|\underline{w}\|} \quad \text{s.t. } \forall i : y_i(\underline{w} \underline{x}_i + b) \geq 1$$

L2 Norm



$$\max_{\underline{w}, b, \gamma} \left[\min_i \frac{|\underline{w} \underline{x}_i + b|}{\|\underline{w}\|} \right]$$

$$\forall i : \frac{|\underline{w} \underline{x}_i + b|}{\|\underline{w}\|} \geq \gamma . \quad \gamma \geq 0$$

$$y_i(\underline{w} \underline{x}_i + b) = |\underline{w} \underline{x}_i + b|$$

$$\forall i: \frac{y_i(\underline{w} \cdot \underline{x}_i + b)}{\|\underline{w}\|} \geq \gamma \rightarrow \forall i: y_i(\underline{w} \cdot \underline{x}_i + b) \geq \gamma \|\underline{w}\|$$

\downarrow norm of \underline{w}

scale $\|\underline{w}\| : \|\underline{w}\| \rightarrow \frac{1}{\gamma}$

\downarrow

$$\max_{\underline{w}, b} \frac{1}{\|\underline{w}\|} \text{ s.t. } \forall i: y_i(\underline{w} \cdot \underline{x}_i + b) \geq 1.$$

find \underline{w}, b

primal SVM.

$$\nexists \left\{ \begin{array}{l} \min_{\underline{w}, b} \frac{1}{2} \|\underline{w}\|^2 \text{ s.t. } \forall i: y_i(\underline{w} \cdot \underline{x}_i + b) \geq 1 \\ m \geq \frac{1}{\varepsilon^2} \min \left[\frac{4}{(\gamma^*)^2} \log^2 \left(\frac{4}{\varepsilon^2 (\gamma^*)^2} \right) d \right] \end{array} \right. \rightarrow \begin{array}{l} \text{solve this and we got} \\ \text{support vectors} \\ y_i(\underline{w} \cdot \underline{x}_i + b) = 1 \end{array}$$

number of data size.

$$\underline{w}^* = \sum_{i=1}^m \alpha_i y_i \underline{x}_i \quad \alpha_i \geq 0$$

$\alpha_i = 0$ iff (x_i, y_i) not a support vector.

$$\begin{aligned} f(x) &= \operatorname{sgn}(\underline{w} \cdot \underline{x} + b) \\ &= \operatorname{sgn} \left(\sum_{i=1}^m \alpha_i y_i \underline{x}_i \cdot \underline{x} + b \right) \\ &= \operatorname{sgn} \left(\sum_{i \in \text{S.V.}} \alpha_i y_i (\underline{x}_i \cdot \underline{x}) + b \right) \quad \text{we only need to consider the support vectors.} \end{aligned}$$

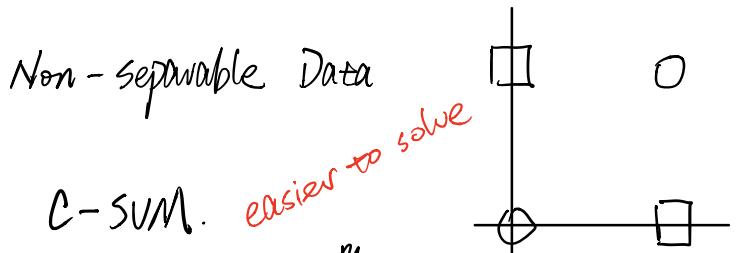
Dual SVM Problem

$$\max_{\underline{\alpha}} \left[\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{j=1}^m \sum_{i=1}^m \alpha_i \alpha_j y_i (\underline{x}_i \cdot \underline{x}_j) y_j \right]$$

$$\text{st. } \sum_{i=1}^m \alpha_i y_i = 0 \quad \forall i: \alpha_i \geq 0.$$

- solve for the dual α_i 's
- Generate $\underline{w}^* = \sum_i \alpha_i y_i \underline{x}_i$
- Pick some i s.t. $\alpha_i > 0$

$$\begin{aligned} y_i (\underline{w}^* \cdot \underline{x}_i + b) &= 1 \\ y_i y_i (\underline{w}^* \cdot \underline{x}_i + b) &= y_i \\ \underline{w}^* \cdot \underline{x}_i + b &= y_i \\ b &= y_i - \underline{w}^* \cdot \underline{x}_i \end{aligned}$$



C-SVM. easier to solve

$$\min_{\underline{w}, b, \underline{\xi}} \frac{1}{2} \|\underline{w}\|^2 + C \left(\frac{1}{m} \sum_{i=1}^m \underline{\xi}_i \right)$$

s.t. $\forall i. y_i (\underline{w} \cdot \underline{x}_i + b) \geq 1 - \underline{\xi}_i$ $\underline{\xi}_i = \text{error}$
 $\forall i. \underline{\xi}_i \geq 0$ how much error we allow.

V-SVM. better result.

$$\begin{array}{ll} \min_{\underline{w}, b, \underline{\xi}, \rho} & \frac{1}{2} \|\underline{w}\|^2 - \rho \rho + \frac{1}{m} \sum_{i=1}^m \underline{\xi}_i & \text{err}_{\text{train}}(f) \leq \rho \\ \text{s.t. } & \forall i. y_i (\underline{w} \cdot \underline{x}_i + b) \geq \rho - \underline{\xi}_i \\ & \forall i. \underline{\xi}_i \geq 0 \\ & \rho \geq 0 \end{array}$$

$$\min \frac{1}{2} \|\underline{w}\|^2 \quad \text{s.t. } 0 \geq 1 - y_i (\underline{w} \cdot \underline{x}_i + b)$$

$$\nabla f(\underline{w}) + \sum_{i=1}^m \lambda_i \nabla g_i(\underline{w}, b) = 0$$

1

$$\begin{aligned} \underline{w} + \sum_{i=1}^m \lambda_i (-y_i x_i) &= 0 \\ 0 + \sum_{i=1}^m \lambda_i (-y_i) &= 0 \\ \underline{w} = \sum_{i=1}^m \lambda_i y_i \underline{x}_i \\ \sum_{i=1}^m \lambda_i y_i &= 0 \end{aligned}$$

Lagrangian Duality

$$\begin{aligned} &\min_{\underline{w}, b} f(\underline{w}, b) \\ \text{s.t. } &\forall i: 0 \geq g_i(\underline{w}, b) \\ \max_{\lambda} &\min_{\underline{w}, b} [f(\underline{w}, b) + \sum_{i=1}^m \lambda_i g_i(\underline{w}, b)] \\ &\lambda_i \geq 0 \end{aligned}$$

Non-linear SVMs

Feature Embedding

$$\phi(x_1, x_2) = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$$

$$x_1 \mid x_2$$

↓

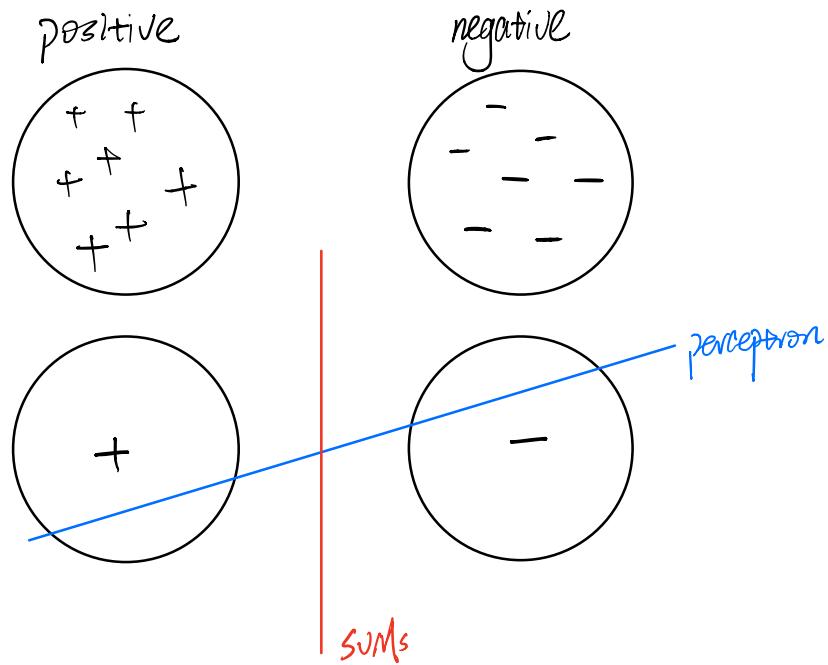
$$x_1 \mid x_2 \mid x_1 \cdot x_2 \mid x_1^2 \mid x_2^2$$

XOR			
-	-	+	-
-	+	+	-
-	-	-	-
+c	-	D	
A: (1, 1)			
B: (-1, 1)			
C: (-1, -1)			
D: (1, -1)			
$\text{sgn}(-x_1 \cdot x_2)$			
$x_1^2 + x_2^2 \leq 1$			
$0 \leq -x_1^2 - x_2^2$			
$\text{sgn}(1 - x_1^2 - x_2^2)$			

~

Sample Complexity

S.C. for perceptions \geq S.C. for SVMs



Optimization

$\min_{\underline{x} \in C} F(\underline{x}) \rightarrow$ Given F.C. find $\min_{\underline{x} \in C} F(\underline{x})$.

Focus on d.

$$C = \mathbb{R}$$

F: differentiable.

Necessary Condition:

Any minimizer must satisfy $\nabla F(\underline{x}^*) = 0$

$$\begin{aligned} F(\underline{x}) &= \|\underline{x} - \underline{a}\|^2 + \|\underline{x} - \underline{b}\|^2 \\ &= (\underline{x} - \underline{a})^T (\underline{x} - \underline{a}) + (\underline{x} - \underline{b})^T (\underline{x} - \underline{b}) \\ &= 2\underline{x}^T \underline{x} - 2\underline{x}^T \underline{a} - 2\underline{x}^T \underline{b} + \|\underline{b}\|^2 + \|\underline{a}\|^2 \end{aligned}$$

$$\nabla \bar{F}(\underline{x}) = 2(\underline{x} - \underline{a}) + 2(\underline{x} - \underline{b}) = 0$$

$$\underline{x} = \frac{1}{2}(\underline{a} + \underline{b})$$

Guess \underline{x} , try to improve it to \underline{x}' where $\bar{F}(\underline{x}') < \bar{F}(\underline{x})$

$$\underline{x}' = \underline{x} + \underline{d}$$

What must be true about \underline{d} so that $\bar{F}(\underline{x}') < \bar{F}(\underline{x})$

$$\underline{x}_1 = \underline{x}_0 + \underline{d}_0$$

$$f(\underline{x}_1) = f(\underline{x}_0) + f'(\underline{x}_0)(\underline{d}_0)$$

Optimization

$$\min_{\underline{x} \in \mathbb{R}^d} \bar{F}(\underline{x})$$

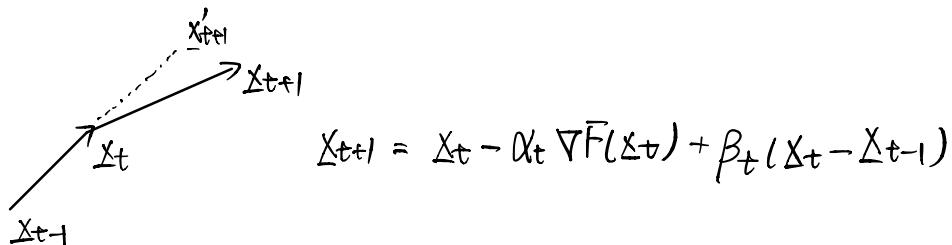
\underline{x}_0 = Some guess

$$\underline{x}_{t+1} = \underline{x}_t - \alpha_t \nabla \bar{F}(\underline{x}_t) \quad \text{gradient direction}$$

$$\underline{x}_t \rightarrow \underline{x}^* \quad \nabla \bar{F}(\underline{x}^*) = 0$$

Coordinate Descent : one coordinate at a time.

Momentum method



$$\bar{F}(\underline{x}) = \sum_{i=1}^m F_i(\underline{x})$$

Stochastic Gradient Descent

$$\bar{F}(\underline{x}) \approx \mathbb{E}[\bar{F}_i(\underline{x})]$$

$$\bar{F}(\underline{x}) = \frac{1}{m} \sum_{i=1}^m \bar{F}_i(\underline{x}) \quad \nabla \bar{F}(\underline{x}) \approx \mathbb{E}[\nabla \bar{F}_i(\underline{x})]$$

SGD.

$\underline{x}_0 = \text{Some Guess}$

At time t , pick i at random

$$\underline{x}_{t+1} = \underline{x}_t - \alpha_t \nabla \bar{F}_i(\underline{x}_t)$$

Total Error

$$\sum_{i=1}^m ((\underline{w} \cdot \underline{x}_i + b) - y_i)^2$$

Constrained Optimization

$$\min_{\underline{x} \in C} \bar{F}(\underline{x})$$

$$\begin{aligned} \min & \frac{1}{2} \|\underline{w}\|^2 \\ \text{s.t. } & \forall i: y_i (\underline{w} \cdot \underline{x}_i + b) \geq 1 \end{aligned}$$



Computing SVMs

$$\max_{\alpha_1, \dots, \alpha_m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i (\underline{x}_i \cdot \underline{x}_j) y_j \alpha_j$$

$$\text{s.t. } \sum_{i=1}^m \alpha_i y_i = 0, \quad \forall i: \alpha_i \geq 0$$

↓

$$\max_{\alpha_1, \dots, \alpha_m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i (\underline{x}_i \cdot \underline{x}_j) y_j \alpha_j$$

$$\text{s.t. } \sum_{\substack{\text{positive} \\ \text{class}}} \alpha_i = \sum_{\substack{\text{negative} \\ \text{class}}} \alpha_j, \quad \forall i: \alpha_i \geq 0$$

SMO

- α
- Pick α_i, α_j
- Modif. together to maintain positive equality
- improve objective function

$$\max_{\alpha_i, \alpha_j} [A + \alpha_i + \alpha_j]$$

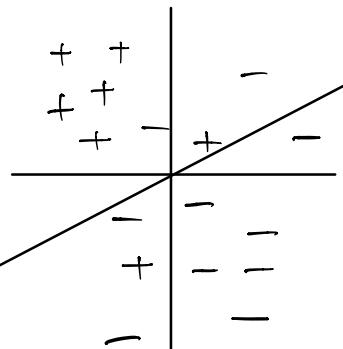
Logistic Regression

$F(x)$ = Probability that x belongs to class +1

$$F(\underline{x}) = \frac{1}{1 + e^{-(w\underline{x} + b)}}$$

$$\text{Sigmoid}(z) = \frac{1}{1 + e^{-z}}$$

$$\{(x_i, y_i)\}, y_i \in \{0, 1\}.$$



$$\begin{aligned} \text{Error}_S(w, b) \\ = \frac{1}{m} \sum_{i=1}^m (y_i - F(\underline{x}_i))^2 \end{aligned}$$

$$\text{lik}(y_1, \dots, y_N | x_1, \dots, x_N) = \prod_{i=1}^N F(\underline{x}_i)^{y_i} (1 - F(\underline{x}_i))^{1-y_i}$$

↓

$$\begin{aligned} \ln \text{lik}(y_1, \dots, y_N | x_1, \dots, x_N) &= \sum_{i=1}^N \ln [F(\underline{x}_i)^{y_i} (1 - F(\underline{x}_i))^{1-y_i}] \\ &\min_{w, b} \frac{1}{N} \sum_{i=1}^N [-y_i \ln F(\underline{x}_i) + (1-y_i) \ln (1 - F(\underline{x}_i))] \end{aligned}$$

$$\text{logistic Error} = -y_i \ln F(\underline{x}_i) - (1-y_i) \ln (1-F(\underline{x}_i))$$

Logistic Regression

$$F_w(\underline{w}) = \sigma(\underline{w} \cdot \underline{x})$$

$$\bar{E}_{\underline{w}}(\underline{w}) = \frac{1}{m} \sum_{i=1}^m [-y_i \ln F_w(\underline{x}_i) - (1-y_i) \ln (1-F_w(\underline{x}_i))]$$

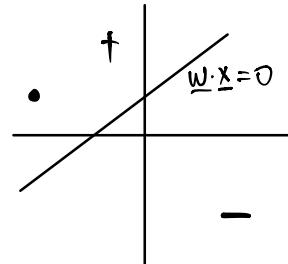
$$\min_{\underline{w}} \bar{E}_{\underline{w}}(\underline{w})$$

$\underline{w}_0 = \text{Some guess}$

At time t , pick $\underline{x}_i \cdot y_i$

$$\underline{w}_{t+1} = \underline{w}_t - \alpha_t \nabla_{\underline{w}} \bar{E}_{\underline{w}}(\underline{w}_t)$$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$



$$\begin{aligned} \frac{\partial \bar{E}_{\underline{w}_i}}{\partial w_j} &= -y_i \cdot \frac{\frac{\partial F(\underline{x}_i)}{\partial w_j}}{F(\underline{x}_i)} - (1-y_i) \left[\frac{-\frac{\partial F(\underline{x}_i)}{\partial w_j}}{1-F(\underline{x}_i)} \right] \\ &= \left[-y_i \frac{1}{F(\underline{x}_i)} + (1-y_i) \frac{1}{1-F(\underline{x}_i)} \right] \underbrace{\frac{\partial F(\underline{x}_i)}{\partial w_j}} \end{aligned}$$

$$\begin{aligned} \frac{\partial F(\underline{x})}{\partial w_j} &= \sigma'(\underline{w}, \underline{x}) \frac{\partial}{\partial w_j} [\underline{w}, \underline{x}] \\ &= \sigma'(\underline{w}, \underline{x}) \cdot x_j \end{aligned}$$

$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$

$$\begin{aligned} \sigma'(\underline{w}, \underline{x}) &= \sigma(\underline{w}, \underline{x})(1-\sigma(\underline{w}, \underline{x})) \\ &= F(\underline{x})(1-F(\underline{x})) \end{aligned}$$

$$\nabla \bar{E}_{\underline{w}_i} = [F(\underline{x}_i) - y_i] \underline{x}_i$$

SGD on Logistic Reg

w_0 : Some guess

At the time t , pick (x_i, y_i)

$$\underline{w}_{t+1} = \underline{w}_t - \alpha_t [F(x_i) - y_i] \underline{x}_i$$

$$F(\underline{x}) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2 + w_{11} x_1^2 + w_{12} x_1 x_2 + w_{22} x_2^2)}}$$

Multiclass Problems

$$\{(x_i, y_i)\} \quad y \in \{1, 2, \dots, C\}$$

probability of \bar{F} in class 1.

$$y=3$$

$$\bar{F}(x) = (\bar{F}_1(x), \bar{F}_2(x), \dots, \bar{F}_C(x))$$

$$y = (0, 0, 1, 0, 0, 0, 0)$$

Cross Entropy Loss

$$Err(\bar{F}) = \frac{1}{m} \sum_{i=1}^m \left[-\sum_{c=1}^C y_i \ln \bar{F}_c(x_i) \right]$$

$$-y_i \ln \bar{F}_c(x_i) - (1-y_i) \ln (1-\bar{F}_c(x_i))$$

$$\{w_1, w_2, \dots, w_C\}$$

$$\bar{F}_c(x) \geq 0$$

$$\bar{F}_c(x) = \frac{e^{w_c \cdot x}}{\sum_{c=1}^C e^{w_c \cdot x}}$$

$$\sum_c \bar{F}_c(x) = 1$$

$$\{0.1, 0.7, 0.93 \quad \bar{F}_c(x) = \frac{1}{1 + e^{-w_c \cdot x}} \quad \sum_c \bar{F}_c(x) = 1}$$

can in multiple class.

$$C=2. \quad \bar{F}_{x_1} = \begin{bmatrix} e^{w_1 \cdot x} & e^{w_2 \cdot x} & 1 \end{bmatrix}$$

$$\left(\frac{e^{\underline{w}_1 \cdot \underline{x}}}{e^{\underline{w}_1 \cdot \underline{x}} + e^{\underline{w}_2 \cdot \underline{x}}}, \frac{e^{\underline{w}_2 \cdot \underline{x}}}{e^{\underline{w}_1 \cdot \underline{x}} + e^{\underline{w}_2 \cdot \underline{x}}} \right)$$

$$= \left(\frac{1}{1 + e^{(\underline{w}_2 - \underline{w}_1) \cdot \underline{x}}}, 1 - \frac{1}{1 + e^{(\underline{w}_2 - \underline{w}_1) \cdot \underline{x}}} \right)$$

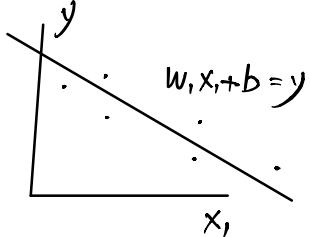
Linear Regression.

$\{(\underline{x}_i, y_i)\}_{i=1}^m \quad y_i \in \mathbb{R}$. Find \underline{w} .

st. $\underline{w} \cdot \underline{x}_i$ matches y_i .

$$\text{err}(\underline{w}) = \frac{1}{m} \sum_{i=1}^m (\underline{w} \cdot \underline{x}_i - y_i)^2$$

$$\min_{\underline{w}} \frac{1}{m} \sum_{i=1}^m (\underline{w} \cdot \underline{x}_i - y_i)^2$$



linear regression \rightarrow regression problem

logistic regression \rightarrow classification problem

Linear Regression

$$\{(\underline{x}_i, y_i)\}_{i=1}^m \quad f(\underline{x}) = \underline{w} \cdot \underline{x}$$

$$\begin{aligned} \text{err}_{\text{train}}(\underline{w}) &= \sum_{i=1}^m (\underline{w} \cdot \underline{x}_i - y_i)^2 \Rightarrow \min_{\underline{w}} \|\underline{X}\underline{w} - \underline{y}\|^2 \\ &= \|\underline{X}\underline{w} - \underline{y}\|^2 \end{aligned}$$

$$F(\underline{w}) = \|\underline{X}\underline{w} - \underline{y}\|^2 = (\underline{X}\underline{w} - \underline{y})^\top (\underline{X}\underline{w} - \underline{y})$$

$$\nabla F(\underline{w}) = 0 = (\underline{w}^\top \underline{X}^\top - \underline{y}^\top)(\underline{X}\underline{w} - \underline{y})$$

$$= \underline{w}^\top \underline{X}^\top \underline{X}\underline{w} - \underline{w}^\top \underline{X}^\top \underline{y} - \underline{y}^\top \underline{X}\underline{w} + \|\underline{y}\|^2$$

$$\begin{aligned}
 &= \underline{w}^T \underline{x}^T \underline{x} \underline{w} - 2 \underline{w}^T \underline{x}^T \underline{y} + \|\underline{y}\|^2 \\
 \nabla \bar{F} &= 2 \underline{x}^T \underline{x} \underline{w} - 2 \underline{x}^T \underline{y} = 0 \\
 \underline{x}^T \underline{x} \underline{w} &= \underline{x}^T \underline{y}
 \end{aligned}$$

$$\underline{w}^* = [\underline{x}^T \underline{x}]^{-1} \underline{x}^T \underline{y}$$

small
noise vector

$$\underline{y} = \underline{x} \underline{w} + \underline{\varepsilon}$$

$$\underline{w}^* = [\underline{x}^T \underline{x}]^{-1} \underline{x}^T (\underline{x} \underline{w} + \underline{\varepsilon})$$

$$\downarrow = [\underline{x}^T \underline{x}]^{-1} (\underline{x}^T \underline{x} \underline{w} + \underline{x}^T \underline{\varepsilon})$$

final model

$$= \underline{w} + [\underline{x}^T \underline{x}]^{-1} \underline{x}^T \underline{\varepsilon}$$

with noise. Small noise in model.

if $\underline{y} = \underline{x} \underline{w} + \underline{\varepsilon}$

$$\underline{w}^* = \underline{w} + [\underline{x}^T \underline{x}]^{-1} \underline{x}^T \underline{\varepsilon} \quad \underline{\varepsilon} \sim N(\underline{0}, \sigma^2 \underline{I})$$

$$\underline{w} \underline{\varepsilon} \sim N(\underline{0}, \underline{w} \underline{w}^T \sigma^2 \underline{I})$$

$$\underline{\varepsilon} \sim N(0, 1)$$

$$m \underline{\varepsilon} \sim N(0, m^2)$$

$$\begin{aligned} \text{Cov}(x, y) \\ = \bar{E} [(x - \bar{m}_x)(y - \bar{m}_y)] \end{aligned}$$

$$\underline{w}^* \sim \underline{w} + [\underline{x}^T \underline{x}]^{-1} \underline{x}^T N(0, \sigma^2 \underline{I}) \quad (x_1, \dots, x_N)$$

$$\sim \underline{w} + N(\underline{0}, \sigma^2 [\underline{x}^T \underline{x}]^{-1} \underline{x}^T \underline{x} [\underline{x}^T \underline{x}]) \quad C: \text{Cov matrix}$$

$$\underline{w}^* \sim \underline{w} + N(\underline{0}, [\underline{x}^T \underline{x}]^{-1} \sigma^2) \quad C_{ij} = \text{Cov}(x_i, x_j)$$

$$\underline{\varepsilon} \sim N(\underline{0}, \sigma^2 \underline{I})$$

$$\Sigma = X^T X \quad \text{symmetric matrix} \rightarrow \Sigma = Q D Q^T$$

$w^* \sim \underline{w} + N(\underline{0}, \Sigma^{-1} \sigma^2)$

Q : eigen vectors of Σ
 D : diagonal matrix of eigen values of Σ .
 $\Sigma^{-1} = Q D^{-1} Q^T$

$$w^* \sim \underline{w} + N(\underline{0}, Q D^{-1} Q^T \sigma^2) \rightarrow \Sigma \sim N(0, 1)$$

$$w^* \sim \underline{w} + \sigma Q N(\underline{0}, D^{-1}) \rightarrow \sigma \Sigma \sim N(0, \sigma^2)$$

$$D^{-1} = \begin{bmatrix} \lambda_1 & & \xrightarrow{\substack{\text{Variance in first \\ component}}} & \Sigma = Q D Q^T \\ \lambda_2 & \ddots & & \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \\ 0 & \dots & \lambda_k & \end{bmatrix}$$

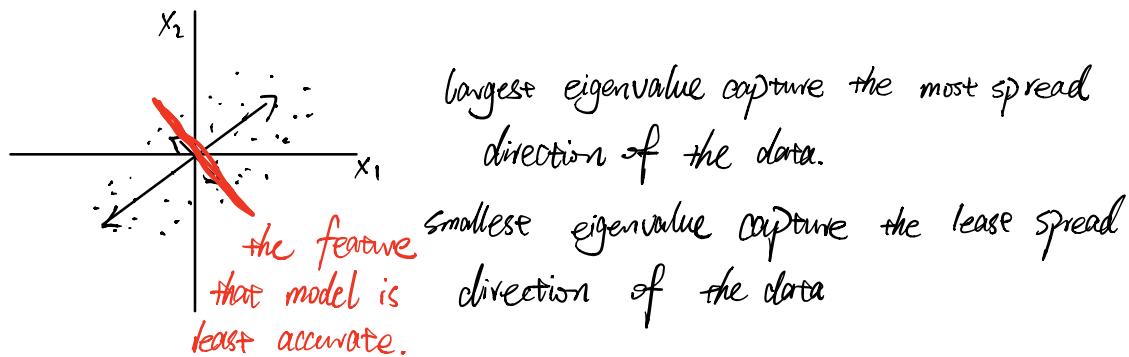
The smaller the eigenvalue of Σ ,
the larger the noise in w^* .

λ_1 very small, λ_k very large.

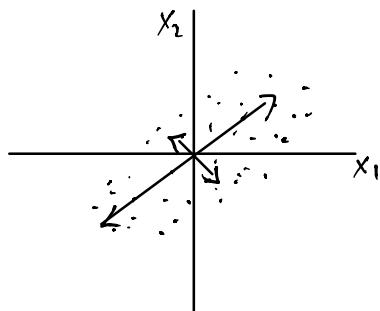
$$K(\Sigma) = \frac{\text{largest eigenvalue}}{\text{smallest eigenvalue.}}$$

If any λ_i is very small, huge error on w^* .

What does it mean if eigenvalues of $\Sigma = X^T X$ are small?



The geometry of data is directly related to the quality of model.



find the largest projection direction of the data.

Find direction \underline{v} , to maximize

$$\sum_{i=1}^m (\underline{v} \cdot \underline{x}_i)^2$$

$$\|\underline{x}\underline{v}\|^2$$

$$\sqrt{\underline{v}^T \underline{x}^T \underline{x} \underline{v}}$$

$$\max_{\underline{v}} \underline{v}^T \underline{x}^T \underline{x} \underline{v}$$

$$\text{s.t. } \|\underline{v}\| = 1$$

Linear Regression

$$\min_{\underline{w}} \sum_{i=1}^m (\underline{w} \cdot \underline{x}_i - y_i)^2$$

$$\min_{\underline{w}} \|\underline{x}\underline{w} - \underline{y}\|^2$$

$$\underline{w}^* = [\underline{x}^T \underline{x}]^{-1} \underline{x}^T \underline{y}$$

$$y = \underline{x}\underline{w} + \underline{\epsilon} \quad \underline{\epsilon} \sim N(0, \sigma^2 I)$$

$$\underline{w}^* \sim \underline{w} + \sigma Q N(0, D^{-1})$$

$$\Sigma = \underline{x}^T \underline{x}$$

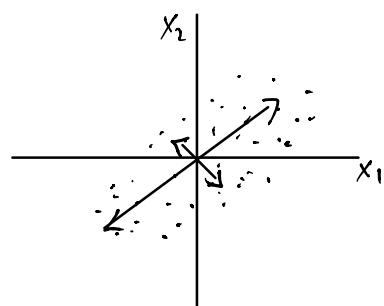
Q : eigenvector of Σ .

D : eigenvalues of Σ .

$$D^{-1} = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_K \end{bmatrix}$$

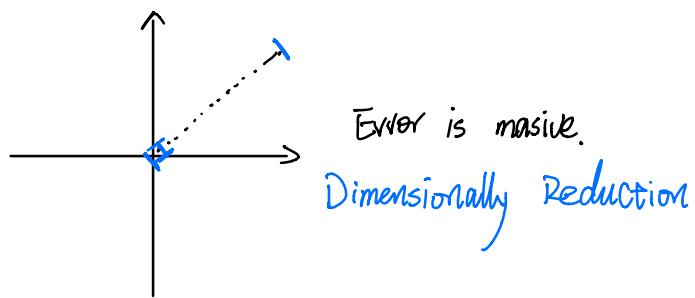
$$\frac{\text{largest error}}{\text{smallest error}} = \frac{\lambda_K}{\lambda_1} = \frac{\lambda_1}{\lambda_K}$$

$$k(\Sigma) = \frac{\lambda_1}{\lambda_K}$$



Sample complexity

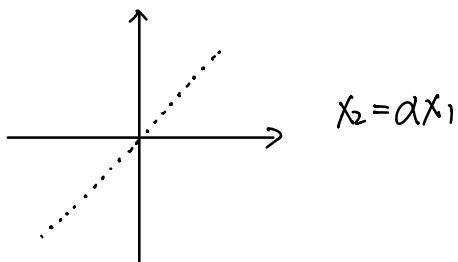
$$m \geq O\left(\frac{H(\varepsilon)}{\varepsilon^2}\right) \Rightarrow \text{err}(\underline{w}^*) \leq \varepsilon^2$$



$$\Sigma = \underline{X}^T \underline{X}$$

$$\Sigma \underline{w} = \underline{X}^T \underline{y}$$

invertible.



$$\min_{\underline{w}} \sum_{i=1}^m (\underline{w}^T \underline{x}_i + y_i)^2 + \lambda \|\underline{w}\|_2^2$$

$$\underline{w}_\lambda^* = [\underline{X}^T \underline{X} + \lambda I]^{-1} \underline{X}^T \underline{y}$$

Linear Regression

$$\min_{\underline{w}} \|\underline{X}\underline{w} - \underline{y}\|_2^2$$

$$\underline{w}^* = [\underline{X}^\top \underline{X}]^{-1} \underline{X}^\top \underline{y}$$

Ridge: $\min_{\underline{w}} \|\underline{X}\underline{w} - \underline{y}\|_2^2 + \frac{\lambda \|\underline{w}\|_2^2}{\text{regulation param}}$

$\underline{X} \rightarrow \text{mean centered}$

$\underline{y} \rightarrow \text{mean centered}$

Feature Selection

$$\underline{y} = \sim + \underbrace{10^{-24} X_5}_{\text{is not important for our model.}} + \sim$$

X_5 is not important for our model.

Pruning - Cut out data.

$$S \subseteq \{1, \dots, k\}.$$

$$\min_{\underline{w}_S} \|\underline{X}_S \underline{w}_S - \underline{y}\|_2^2 \quad \text{Find } S \text{ to minimize this.}$$

subsets are too much!

Lasso: $\min_{\underline{w}} \|\underline{X}\underline{w} - \underline{y}\|_2^2 + \lambda \|\underline{w}\|_1$ Automatically pruned.

$\|\underline{w}\|_1 = \sum |w_i|$

$$\min_{\underline{w}} (\underline{w} - \alpha)^2 + \lambda |\underline{w}|$$

$$\underline{w}^* = \begin{cases} \alpha - \frac{\lambda}{2} & \alpha > \frac{\lambda}{2} \\ \alpha + \frac{\lambda}{2} & \alpha < -\frac{\lambda}{2} \\ 0 & \alpha \in [-\frac{\lambda}{2}, \frac{\lambda}{2}] \end{cases}$$

$$\underline{w} = (1, 0, 0, 0, \dots, 0) \quad \leftarrow \text{sparse vector}$$

$$\underline{v} = \left(\frac{1}{\sqrt{N}}, \frac{1}{\sqrt{N}}, \dots, \frac{1}{\sqrt{N}} \right) \quad \leftarrow \text{even vector.}$$

$$\|\underline{w}\|_2^2 = 1 \quad \|\underline{w}\|_1 = 1$$

$$\|\underline{v}\|_2^2 = 1 \quad \|\underline{v}\|_1 = \sqrt{N}$$

$$\text{Lasso: } \min_{\underline{w}} \|\underline{X}\underline{w} - \underline{y}\|^2 + \lambda \|\underline{w}\|_1$$

support: x_i that actually matter to y .

conditions for support recovery

$$y = \underline{X}_S \cdot \underline{w}_S + \underline{\epsilon}$$

\uparrow
random noise

- Mutual Incoherence

- Minimum Spread: $\lambda_{\min}(X_S^T X_S)$ sufficient large.

- Minimum signal: For any $j \in S^*$ $|w_j|$ can't be too small.

Boosting
 $f_1(x), f_2(x) \dots f_k(x)$ why multiple model.
 $f(x) = \text{majority}(f_1(x), f_2(x) \dots f_k(x))$ different hypothesis space

Ada boost → mega model. train new model from old ones.
adjustment over old model.

$$F_T(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t f_t(x)\right) \quad Q_T(x) = \sum_{t=1}^T \alpha_t f_t(x)$$

$\text{sign } \alpha_1, \dots, \alpha_T, f_1 \dots f_T$

find α_{T+1}, f_{T+1}

$$F_{T+1}(x) = \text{sign}\left(\sum_{t=1}^{T+1} \alpha_t f_t(x)\right)$$



$$\bar{F}_T(x) = \text{sign}(Q_T(x))$$

Erron func

for binary classification

is discrete.

$$\text{err}_{\text{train}}(f) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}\{f(x_i) \neq y_i\}.$$

$$E_T = \sum_{(x,y)} e^{-y Q_T(x)} \quad \rightarrow \begin{array}{l} \text{correct: } E_T \geq 1 \\ \text{incorrect: } E_T \leq 1 \end{array}$$

$$\bar{E}_{T+1} = \sum_{(x,y)} e^{-y(Q_T(x) + \alpha_{T+1} f_{T+1}(x))}$$

$$= \sum_{(x,y)} e^{-y Q_T(x)} e^{-\alpha_{T+1} f_{T+1}(x)} y$$

$$= \left(\sum_{f_{T+1}(x)=y} e^{-y Q_T(x)} \right) e^{-\alpha_{T+1}} + \left(\sum_{f_{T+1}(x) \neq y} e^{-y Q_T(x)} \right) e^{\alpha_{T+1}}$$

what α_{T+1} minimize \bar{E}_{T+1}

$$G_{T+1} = \sum_{f_{T+1}(x) = y} e^{-y Q_T(x)}$$

$$\bar{E}_{T+1} = G_{T+1} e^{-\alpha_{T+1}} + B_{T+1} e^{\alpha_{T+1}}$$

$$B_{T+1} = \sum_{f_{T+1}(x) \neq y} e^{-y Q_T(x)}$$

see derivative 0
and solve

$$\alpha_{T+1} = \frac{1}{2} \ln \left(\frac{G_{T+1}}{B_{T+1}} \right)$$

$$\bar{E}_{T+1} = 2 \sqrt{G_{T+1} B_{T+1}}$$

$$\bar{E}_{T+1} = \sum_{(x,y)} e^{-y(Q_T(x) + \alpha_{T+1} f_{T+1}(x))}$$

$$\bar{E}_T = \sum_{(x,y)} e^{-y Q_T(x)}$$

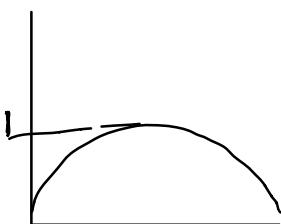
$$= G_{T+1} + B_{T+1}$$

$$\bar{E}_{T+1} = 2 \sqrt{(\bar{E}_T - B_{T+1}) B_{T+1}}$$

$$\bar{E}_{T+1} = 2 \bar{E}_T \sqrt{\left(1 - \frac{B_{T+1}}{\bar{E}_T}\right) \frac{B_{T+1}}{\bar{E}_T}}$$

$$Q_T(x) = \sum_{t=1}^T \alpha_t f_t(x)$$

$$2 \sqrt{(1-p)p}$$



To minimize \bar{E}_{T+1} , find f_{T+1} to minimize $\frac{B_{T+1}}{\bar{E}_T}$.

$$\frac{B_{T+1}}{E_T} = \frac{\sum_{(x,y)} e^{-yQ_T(x)} \cdot \mathbb{I}\{f_{T+1}(x) \neq y\}}{\sum_{(x,y)} e^{-yQ_T(x)}}$$

weighted error function.

$$= \frac{\sum_{(x,y)} w(x,y) \mathbb{I}\{f_{T+1}(x) \neq y\}}{\sum_{(x,y)} w(x,y)} = \text{err}_{w^*}^{\text{train}}(f_{T+1})$$

minimize $w(x,y) \rightarrow e^{-yQ_T(x)}$

Ada Boost

- Initiate $w_1(x,y) = 1$. for all $(x,y) \in \text{Data}$
- At time t , find f_t to minimize $\text{err}_{w^*}^{\text{train}}(f)$
- Set $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \text{err}_{w^*}^{\text{train}}(f_t)}{\text{err}_{w^*}^{\text{train}}(f_t)} \right)$
- $w_{t+1}(x,y) = w_t(x,y) e^{-\alpha_t y f_t(x)} \text{err}_w^{\text{train}}(f) = \frac{\sum_{(x,y)} w(x,y) \mathbb{I}\{f_t(x) \neq y\}}{\sum_{(x,y)} w(x,y)}$
- $G(f_1 \dots f_T)$
- $F(x) = \text{sign}(\sum_{t=1}^T \alpha_t f_t(x))$

As long as $\text{err}_{w^*}^{\text{train}}(f_t) < \frac{1}{2}$, we have improvement.

If $\text{err}_{w^*}^{\text{train}}(f_t) \leq \frac{1}{2} - \delta$. for $t = 1 \dots T$.

$$\text{err}_{\text{train}}(F) \leq e^{-2T\delta^2}$$

$f_1 \dots f_T$ - different hypothesis space
no matter what hypothesis space for new model. random pick?

Boosting

$$f_1, f_2, \dots, f_T$$

$$\bar{F}(\underline{x}) = \text{sign} \left(\sum_{t=1}^T \alpha_t f_t(\underline{x}) \right)$$

Gradient Boosting

$$\bar{E}_{\text{err}} = \frac{m}{n} \sum_{i=1}^n L(y_i, \bar{F}(\underline{x}_i))$$

$$\underline{x} = \underline{4}$$

$$\bar{F}(4)_{\text{raw}} = \bar{F}(4) - \lambda \frac{\partial \bar{E}_{\text{err}}}{\partial \bar{F}(4)}$$

Residuals.

$$r_1 = \frac{\partial \bar{E}_{\text{err}}}{\partial \bar{F}(\underline{x}_1)}$$

$$r_2 = \frac{\partial \bar{E}_{\text{err}}}{\partial \bar{F}(\underline{x}_2)}$$

⋮

$$r_m = \frac{\partial \bar{E}_{\text{err}}}{\partial \bar{F}(\underline{x}_m)}$$

Fit a new model

$$\underline{x}_i \rightarrow r_i$$

$$\tilde{\bar{F}}(\underline{x}_i) = r_i$$

$$\bar{F}_{\text{meta}}(\underline{x}) = \bar{F}(\underline{x}) - \lambda \tilde{\bar{F}}(\underline{x})$$

— Train f to minimize $\sum_i L(y_i, f(\underline{x}_i))$

— Set $\lambda_1 = 1$

— At time $t \geq 1$, $\bar{F}_t = \lambda_1 f_1 + \lambda_2 f_2 + \dots + \lambda_t f_t$

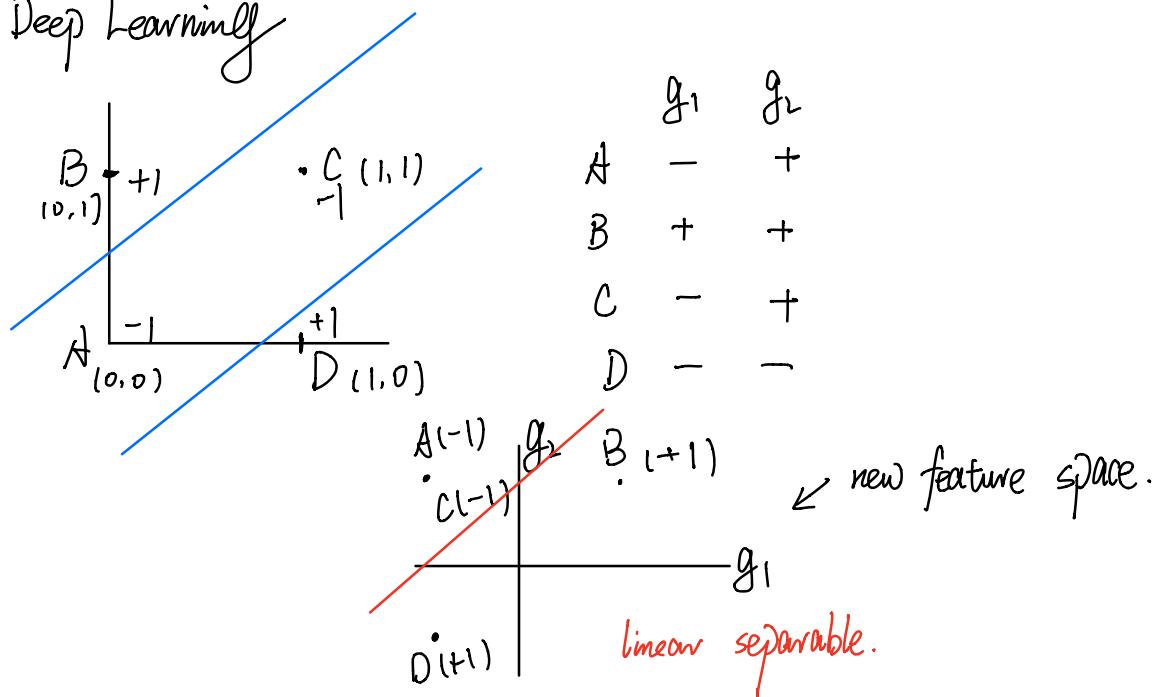
— Define residuals $r_t^i = \frac{\partial}{\partial \bar{F}_t(\underline{x}_i)} [L(y_i, \bar{F}_t(\underline{x}))]$

— Fit f_{t+1} to $\{(x_i, r_t^i)\}$

- Find λ_{t+1} to minimize $\sum_i L(y_i, \bar{F}_t(x_i) - \lambda f_{t+1}(x_i))$
- $\bar{F}_{t+1} = \bar{F}_t - \lambda f_{t+1}$

XGBoost \curvearrowleft apply on decision tree.

Deep Learning



$$g_1(x_1, x_2) = \text{sign}(ax_1 + bx_2 + c)$$

$$g_2(x_1, x_2) = \text{sign}(dx_1 + ex_2 + f)$$

compute features over features

$$\bar{F}(x_1, x_2) = \text{sign}(hg_1(x_1, x_2) + lg_2(x_1, x_2) + p)$$

make computer pick
the features.

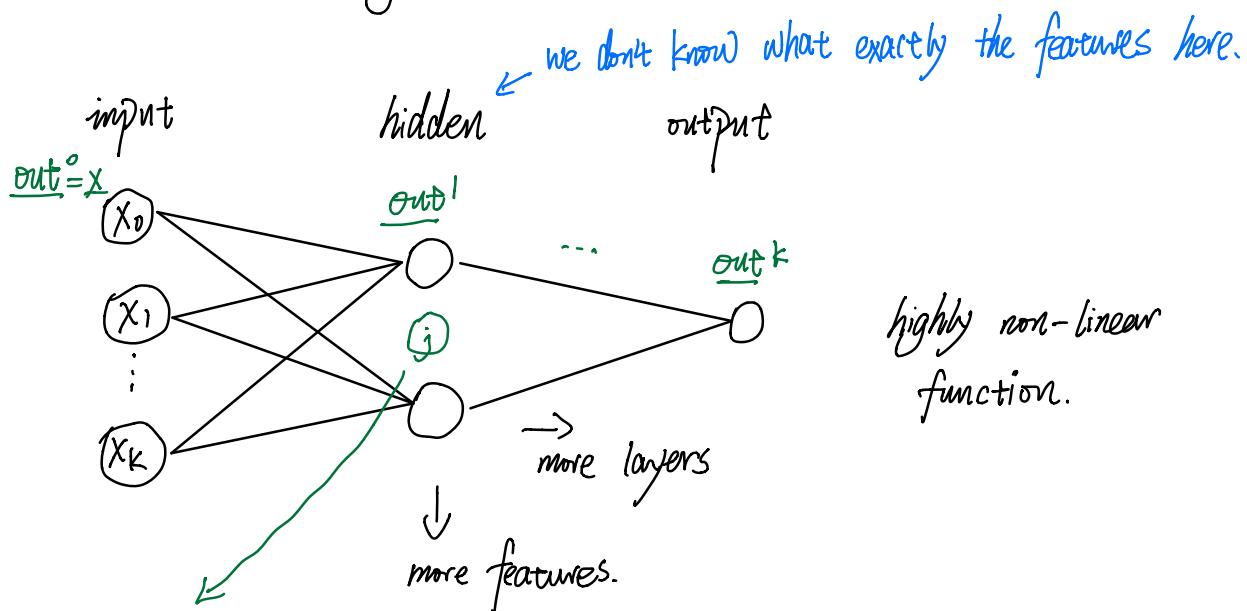
Find $a, b, c, d, e, f, h, l, p \rightarrow$ fit data

NOT find features
ourselves.

Neural Networks

$$f(x_1, x_2) = \sigma(\alpha_1 + \alpha_2 \sigma(\alpha_4 + \alpha_5 x_1 + \alpha_6 x_2) + \alpha_3 \sigma(\alpha_7 + \alpha_8 x_1 + \alpha_9 x_2))$$

$$\sigma(z) = \text{sigmoid}(z)$$



$$\text{out}_{ij}^t = \sigma(w^{t-1}(j) \cdot \text{out}^{t-1})$$

out⁰ = x
↑
dot product.
↓
out^k

Activation Functions.

$$\sigma(z) = \text{sign}(z)$$

$$\sigma(z) = \text{sigmoid}$$

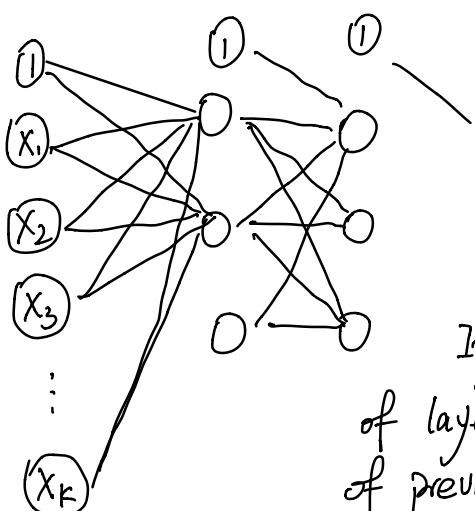
$$= \frac{1}{e^z + 1}$$

$$\sigma(z) = \tanh(z)$$

$$\text{ReLU}$$

$$\sigma(z) = \max(0, z)$$

computational quite good.



If no activation function, backwards of layers are just linear combination of previous layers.

A deeper network allow us compute less features in single layer
 This practically give good result.

Regression

$$\text{Error} = \frac{m}{\sum_{i=1}^m} \|\underline{\text{out}}^k(x_i) - \underline{y}_i\|_2^2$$

Classification.

$$\underline{\text{out}}_c^K = \frac{e^{\underline{\text{out}}_c^{K-1}}}{\sum_c e^{\underline{\text{out}}_c^{K-1}}}$$

$$L = \frac{m}{\sum_{i=1}^m} - \left[\sum_c y_i^c \ln \underline{\text{out}}_c^K(x_i) \right]$$

$$W_{\text{new}} = W_{\text{old}} - \alpha \nabla_{W_{\text{old}}} L$$

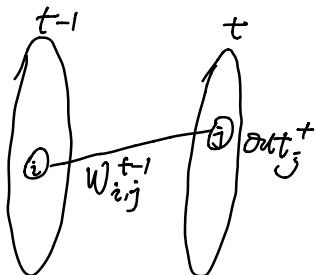
$$\text{Goal: } \frac{\partial L}{\partial w_{i,j}^{t-1}} = \frac{\partial L}{\partial \underline{\text{out}}_j^t} \cdot \frac{\partial \underline{\text{out}}_j^t}{\partial w_{i,j}^{t-1}}$$

$$\underline{\text{out}}_j^t = \sigma(\underline{w}^{t-1}(j) \cdot \underline{\text{out}}^{t-1})$$

$$\frac{\partial \underline{\text{out}}_j^t}{\partial w_{i,j}^{t-1}} = \frac{\partial}{\partial w_{i,j}^{t-1}} [\underline{w}^{t-1}(j) \cdot \underline{\text{out}}^{t-1}] \cdot \sigma(\underline{w}^{t-1}(j) \cdot \underline{\text{out}}^{t-1})$$

$$\frac{\partial \underline{\text{out}}_j^t}{\partial w_{i,j}^{t-1}} = \underline{\text{out}}_i^{t-1} \cdot \sigma(\underline{w}^{t-1}(j) \cdot \underline{\text{out}}^{t-1})$$

$$\frac{\partial L}{\partial w_{i,j}^{t-1}} = \left(\frac{\partial L}{\partial \underline{\text{out}}_j^t} \right) \sigma(\underline{w}^{t-1}(j) \cdot \underline{\text{out}}^{t-1}) \underline{\text{out}}_j^{t-1}$$

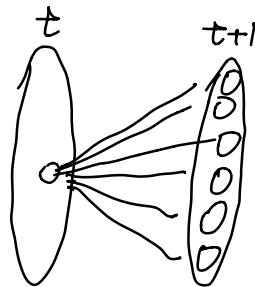


$t = K$: compute directly from L

$t < K$.

$$\frac{\partial L}{\partial \text{out}_j^t} = \sum_{\substack{\text{nodes } k \\ \text{in the next} \\ \text{layer}}} \left(\frac{\partial L}{\partial \text{out}_k^{t+1}} \right) \left(\frac{\partial \text{out}_k^{t+1}}{\partial \text{out}_j^t} \right)$$

backpropagation



$$\text{out}_k^{t+1} = \sigma(\underline{w}_{(k)}^t \cdot \text{out}^t)$$

$$\frac{\partial \text{out}_k^{t+1}}{\partial \text{out}_j^t} = w_{j,k}^t \sigma'(\sim)$$

back-propagation \rightarrow gradient descent. \rightarrow minimize Loss.