

Linear Regression

Input : Vector of **feature** measured by real number

Output : single scalar real number

Example : Predicting weight based on height

Input : **height** of one individual

Output : **weight** of the same individual

Data set

height | weight Goal : Predict weight given height

180	96
162	55
183	96.5
174	70
160	62
163	54
190	60
165	72
175	93
170	99
170	60
169	92
168	59
175	75
169	56
171	99
155	45
158	60
175	60
165	72

Hypothesis : A specific predictor

x : Domain of our hypothesis function (input) \mathbb{R}
(more input e.g. 2 : \mathbb{R}^2)

y : Range of our hypothesis function (output) \mathbb{R}

each hypothesis is a function

h : $\mathbb{R} \rightarrow \mathbb{R}$

What is the form of h ?

$$h(x) = \theta_0 + \theta_1 x_1$$

predicted weight parameter height of h (weights)

(x, y) - one training example

$(x^{(i)}, y^{(i)})$ - i^{th} training example

$$\begin{aligned} x^{(1)} &= 180 & \theta^T &= [\theta_0 \ \theta_1 \dots \ \theta_n] x &= \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \\ x^{(2)} &= 162 \\ y^{(1)} &= 96 \end{aligned}$$

What if we have n inputs ? $n+1$ parameters

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x \quad (\text{dot product})$$

$$= \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$h(x) = \theta_0 + \theta_1 x_1 - \text{Linear equation}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix}$$

Hypothesis : Output(y) is a linear function
of the inputs (x)

How to find the best possible θ ?

Cost function : A function of θ that

squared error function tell us how BAD (GOOD) θ is

Regression problem normally use the

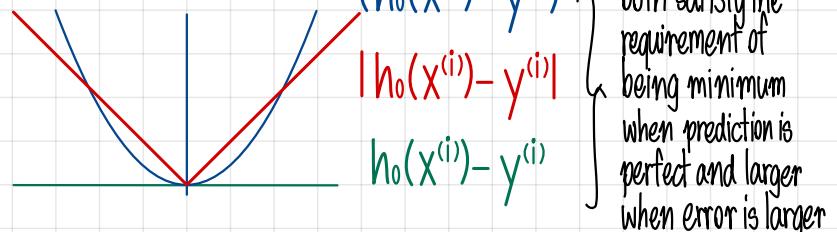
Least Square Cost function (J)

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

our hypothesis function with specific value for θ input feature of i^{th} example observed target of i^{th} example

$$h(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

minimize $J(\theta_0, \theta_1, \theta_2, \dots, \theta_m)$



both satisfy the requirement of being minimum when prediction is perfect and larger when error is larger

So far, we have

1. Form of the hypothesis function

2. Cost function

3. Form of dataset

Still need :

1. An algorithm for minimizing the cost function

To find the minimum of J :

$$\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta) \quad \nabla_{\theta} J(\theta)$$

$J(\theta)$ is minimized by $\frac{\partial J}{\partial \theta_0} = 0, \frac{\partial J}{\partial \theta_1} = 0$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial J}{\partial \theta} = 2U \cdot U'$$

$$\frac{\partial J}{\partial \theta_0} = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot \frac{\partial h_{\theta}(x^{(i)})}{\partial \theta_0}$$

$$\frac{\partial J}{\partial \theta_1} = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot \frac{\partial h_{\theta}(x^{(i)})}{\partial \theta_1}$$

$$0 = \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot 1$$

$$0 = \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot X_1^{(i)}$$

Substitute $\hat{y}^{(i)}$ for $h_{\theta}(x^{(i)})$

$$\sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) = 0$$

$$\sum_{i=1}^m (\hat{y}^{(i)} X_1^{(i)} - y^{(i)} X_1^{(i)}) = 0$$

$$\sum_{i=1}^m \hat{y}^{(i)} - \sum_{i=1}^m y^{(i)} = 0$$

$$\sum_{i=1}^m \hat{y}^{(i)} X_1^{(i)} - \sum_{i=1}^m y^{(i)} X_1^{(i)} = 0$$

$$\text{so } \sum_{i=1}^m \hat{y}^{(i)} = \sum_{i=1}^m y^{(i)}, \sum_{i=1}^m \hat{y}^{(i)} X_1^{(i)} = \sum_{i=1}^m y^{(i)} X_1^{(i)}$$

$$\hat{y}^{(i)} = \theta_0 X_0^{(i)} + \theta_1 X_1^{(i)}$$

$$\sum_{i=1}^m \hat{y}^{(i)} = \sum_{i=1}^m y^{(i)} ; \hat{y}^{(i)} = \theta_0 X_0^{(i)} + \theta_1 X_1^{(i)}$$

$$\sum_{i=1}^m \theta_0 X_0^{(i)} X_0^{(i)} + \sum_{i=1}^m \theta_1 X_1^{(i)} X_0^{(i)} = \sum_{i=1}^m y^{(i)} X_0^{(i)}$$

$$\sum_{i=1}^m \hat{y}^{(i)} X_1^{(i)} = \sum_{i=1}^m y^{(i)} X_1^{(i)}$$

$$\sum_{i=1}^m \theta_0 X_0^{(i)} X_1^{(i)} + \sum_{i=1}^m \theta_1 X_1^{(i)} X_1^{(i)} = \sum_{i=1}^m y^{(i)} X_1^{(i)}$$

$$\theta_0 \sum_{i=1}^m X_0^{(i)} X_1^{(i)} + \theta_1 \sum_{i=1}^m X_1^{(i)} X_1^{(i)} = \sum_{i=1}^m y^{(i)} X_1^{(i)}$$

$$\mathbf{x} = \begin{bmatrix} X_0^{(1)} & X_1^{(1)} \\ X_0^{(2)} & X_1^{(2)} \\ X_0^{(3)} & X_1^{(3)} \\ \dots & \dots \\ X_0^{(m)} & X_1^{(m)} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} Y^{(1)} \\ Y^{(2)} \\ Y^{(3)} \\ \dots \\ Y^{(m)} \end{bmatrix}$$

$$\mathbf{x}^T \mathbf{x} = \begin{bmatrix} \sum_{i=1}^m X_0^{(i)} X_0^{(i)} & \sum_{i=1}^m X_0^{(i)} X_1^{(i)} \\ \sum_{i=1}^m X_1^{(i)} X_0^{(i)} & \sum_{i=1}^m X_1^{(i)} X_1^{(i)} \end{bmatrix} \quad \mathbf{x}^T \mathbf{y} = \begin{bmatrix} \sum_{i=1}^m X_0^{(i)} Y^{(i)} \\ \sum_{i=1}^m X_1^{(i)} Y^{(i)} \end{bmatrix}$$

$$\mathbf{x}^T \mathbf{x} \cdot \theta = \mathbf{x}^T \mathbf{y}$$

$$\begin{bmatrix} \sum_{i=1}^m X_0^{(i)} X_0^{(i)} & \sum_{i=1}^m X_0^{(i)} X_1^{(i)} \\ \sum_{i=1}^m X_1^{(i)} X_0^{(i)} & \sum_{i=1}^m X_1^{(i)} X_1^{(i)} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m X_0^{(i)} Y^{(i)} \\ \sum_{i=1}^m X_1^{(i)} Y^{(i)} \end{bmatrix}$$

$$\mathbf{x}^T \mathbf{x} \cdot \theta = \mathbf{x}^T \mathbf{y} \rightarrow \theta = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y} \quad \text{Normal equation}$$

If $\mathbf{x}^T \mathbf{x}$ has a matrix inverse!

$$\mathbf{x} = \begin{bmatrix} 1 & X_1^{(1)} \\ 1 & X_1^{(2)} \\ 1 & X_1^{(3)} \\ \dots & \dots \\ 1 & X_1^{(m)} \end{bmatrix} \quad m \text{ examples, } n=1$$

Exercise : Form matrices \mathbf{x}, \mathbf{y} for our weight-height dataset and find optimal values of θ_0, θ_1 . Before tomorrow!

Last time → sample machine learning problem
(predict weight given height)

the 4 things you need for a machine learning

1. Dataset & its structure (x, y)

2. Form of the hypothesis function $h : X \rightarrow Y$
(parameterized by θ)

3. Cost function

4. An algorithm for minimizing the cost function

4 methods: Analytical approach

Batch Gradient Descent

Stochastic Gradient Descent

Mini-Batch Gradient Descent

use one training example at A time to update θ

faster than batch but noiser

Why squared error? $(h_\theta(x^{(i)}) - y^{(i)})^2$?

→ we only have to minimize

→ smooth & nice

→ analytical optimum $\theta^* = (x^T x)^{-1} x^T y$

→ can be derived from probabilistic point of view

New Cost function (Objective function): Likelihood

Generative model for the provenance of our dataset

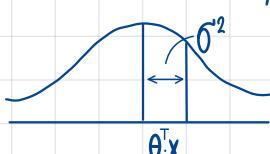
independent & identically distributed

1. sample X i.i.d from some unknown distribution over X

$$X \sim P(X)$$

2. sample y from a gaussian distribution

with mean of $\theta^T x$, variance of σ^2


$$P(Y=y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\theta^T x)^2}{2\sigma^2}}$$

$$y \sim N(\theta^T x, \sigma^2) \text{ normal distribution}$$

standardization of
a normal distribution

Given a dataset, we will measure the goodness or quality of a particular value of θ in terms of how consistent the data are with a gaussian distribution with mean $\theta^T x$ and variance σ^2



Linear Regression

$x : \mathbb{R}$, $y : \mathbb{R}$ after augmenting x with dummy variable 1

$$h(\theta) : \underline{x} \rightarrow \theta^T \underline{x}$$

$$J(\theta : X, y) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta) = (x^T x)^{-1} x^T y$$

Alternative optimization algorithm:

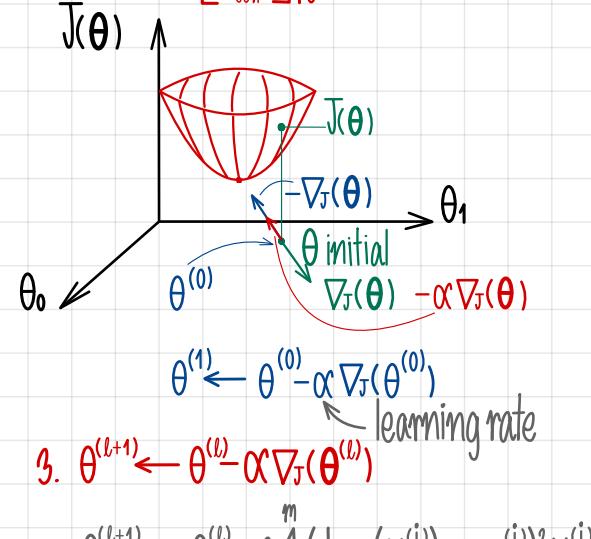
Gradient Descent

1. pick an initial θ

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

2. calculate gradient

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix}_{\theta}$$



3. $\theta^{(l+1)} \leftarrow \theta^{(l)} - \alpha \nabla J(\theta^{(l)})$

$$\theta^{(l+1)} \leftarrow \theta^{(l)} - \alpha \sum_{i=1}^m (h_{\theta^{(l)}}(x^{(i)}) - y^{(i)})^2 x_i^{(i)}$$

4. Repeat from #2 until convergence

Likelihood of θ :

outcomes of a random experiment

$$L(\theta : X, y) = P(X, y ; \theta)$$

parameters

contextual data/ constants

parameters of

the distribution

used in the function

$$\theta^* = \underset{\theta}{\operatorname{argmax}} L(\theta : X, y)$$

$$= \underset{\theta}{\operatorname{argmin}} -L(\theta : X, y)$$

Maximum Likelihood Approach

Choose θ to maximize $L(\theta)$

New Objective function : Likelihood of θ
function of θ parameter

$$\theta^* = \underset{\theta}{\operatorname{argmax}} L(\theta; \mathbf{x}, \mathbf{y})$$

$$L(\theta; \mathbf{X}, \mathbf{y}) = P(\mathbf{X}, \mathbf{y}; \theta) \quad \mathbf{X} \quad \mathbf{y}$$

$$= P(X^{(1)}, X^{(2)}, \dots, X^{(m)}, Y^{(1)}, Y^{(2)}, \dots, Y^{(m)}; \theta)$$

$$P(A, B) = P(A) \cdot P(B) \quad \text{if } A, B \text{ mutually independent}$$

Example: two flips of a fair coin

$$P(H, H) = P(H) P(H) = 0.5 \times 0.5 = 0.25$$

since $(X^{(i)}, Y^{(i)})$, $(X^{(j)}, Y^{(j)})$, $i \neq j$ are sampled

i.i.d. from the same distribution, we have

$$P(X^{(i)}, X^{(j)}, Y^{(i)}, Y^{(j)}) = P(X^{(i)}, Y^{(i)}) \cdot P(X^{(j)}, Y^{(j)})$$

$$P(\mathbf{X}, \mathbf{y}; \theta) = \prod_{i=1}^m P(X^{(i)}, Y^{(i)}; \theta) \quad \text{change form}$$

Are $X^{(i)}$ and $Y^{(i)}$ independent? Nope

$$P(X^{(i)}, Y^{(i)}) \neq P(X^{(i)}) \cdot P(Y^{(i)})$$

$$P(A|B) = P(A, B) \xrightarrow{\substack{\text{factorization of joint probability} \\ P(A)}} P(A, B) = P(A, B) P(A)$$

$$P(X^{(i)}, Y^{(i)}) = P(Y^{(i)}|X^{(i)}) \cdot P(X^{(i)})$$

$$L(\theta; \mathbf{X}, \mathbf{y}) = \left[\prod_{i=1}^m P(Y^{(i)}|X^{(i)}; \theta) \right] \left[\prod_{i=1}^m P(X^{(i)}; \theta) \right]$$

by assumption of our generative model for linear regression
 $Y^{(i)} \sim N(\Theta^T \mathbf{x}^{(i)}, \sigma^2)$

$$P(X^{(i)}|Y^{(i)}; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}}$$

$$L(\theta; \mathbf{X}, \mathbf{y}) = \left[\prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}} \right] \left[\prod_{i=1}^m P(X^{(i)}; \theta) \right]$$

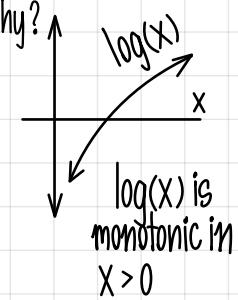
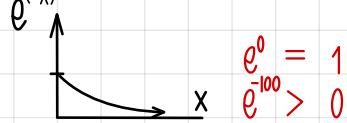
is maximized by θ maximizing

$$\prod_{i=1}^m e^{-\frac{(Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}}$$

is maximized by θ maximizing

$$l(\theta) = \log L(\theta) \log \prod_{i=1}^m e^{-\frac{(Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}}$$

even negative? \times



$$\log(a \cdot b) = \log(a) + \log(b)$$

$$\log \prod_{i=1}^m e^{-\frac{(Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}} = \sum_{i=1}^m \log e^{-\frac{(Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}} = \sum_{i=1}^m -\frac{(Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}$$

$$l(\theta) = -\frac{1}{2\sigma^2} \sum_{i=1}^m (Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2$$

is maximized by θ maximizing

$$-\frac{1}{2} \sum_{i=1}^m (Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2$$

is maximized by θ minimizing

$$-\frac{1}{2} \sum_{i=1}^m (Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2 = \frac{1}{2} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - Y^{(i)})^2$$

$$h_\theta(\mathbf{x}^{(i)}) = J(\theta)$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (Y^{(i)} - \Theta^T \mathbf{x}^{(i)})^2$$

Maximizing the Likelihood of θ under the gaussian

Generative model give us the least square cost

function (J) for Linear Regression !!!

Squared error cost function

1. make sense!

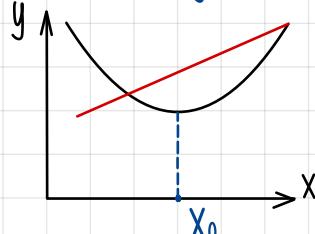
2. give us an analytical solution to the optimization problem (convenient)

3. maximizes the Likelihood of θ

(give us the θ most likely to have been used nature to generate with dataset \mathbf{x}, \mathbf{y})

This approach is called the Generative approach to defining ML algorithms

Polynomial Regression



Quadratic (parabola)
linear

$$y = (x - x_0)^2 + y_0 + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

$$y = \alpha x^2 + \beta x + c + \epsilon$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

linear in parameter
nonlinear in input features

Quadratic Polynomial

in two features: (x_0, x_1)

$$y = \alpha x_0^2 + \beta x_0 + \gamma x_1 + \delta x_0 x_1 + \epsilon$$

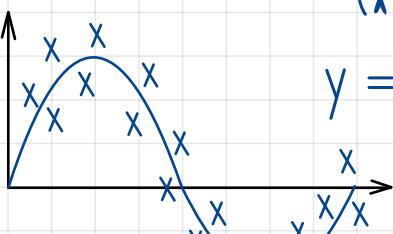
$$\epsilon \sim N(0, \sigma^2)$$

new features are nonlinear combinations of original features

$$X = \begin{bmatrix} 1 & x_0^{(1)} x_0^{(1)} & x_1^{(1)} x_1^{(1)} & x_0^{(1)} x_1^{(1)} & x_0^{(1)} & x_1^{(1)} \\ 1 & x_0^{(m)} x_0^{(m)} & x_1^{(m)} x_1^{(m)} & x_0^{(m)} x_1^{(m)} & x_0^{(m)} & x_1^{(m)} \end{bmatrix}$$

$$Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta) = \underset{\theta}{\operatorname{argmax}} l(\theta) = (X^T X)^{-1} X^T Y$$

$$y = \theta_0 + \theta_1 \cos x + \epsilon$$



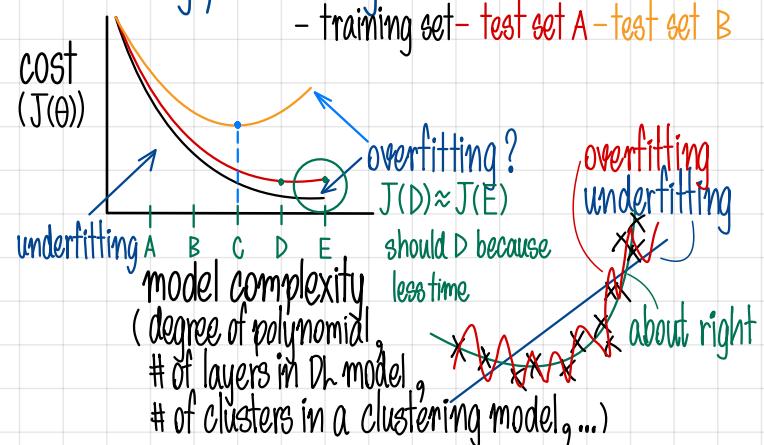
Possible nonlinear relationship between x, y

Q: when to stop adding parameters ???

Overfitting: too many parameters

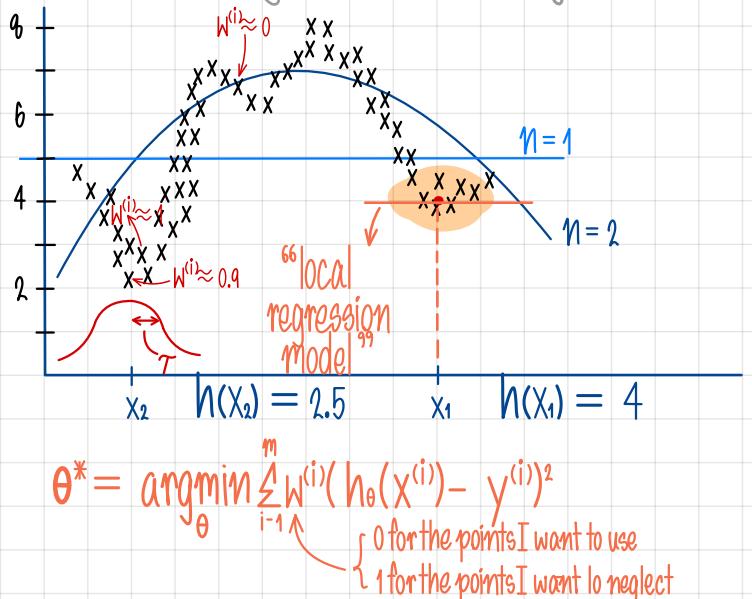
Underfitting: too few parameters

Underfitting, Overfitting



consider

Locally Weighted Linear Regression



Idea: at this time, when asked what $h(x)$ should be for a specific x , we do local linear regression, then answer with the prediction based on that local linear regression result.

How can we perform a linear regression using a subset of the training data?

How can we set the weights $w^{(i)}$ to appropriate values?

"soft" weights are also possible: big for points $x^{(i)}$ close to x and small for far from x

Gaussian kernel for $w^{(i)}$:

$$W^{(i)} = e^{-\frac{(x^{(i)} - x)^T (x^{(i)} - x)}{2T^2}}$$

X: Query point

location x
T "bandwidth" parameter tau.
std. deviation if this was a probability distribution

If $x^{(i)}$ is on top of x (small), $W^{(i)} \approx 1$
If $x^{(i)} - x$ is large (negative or positive directional), $W^{(i)} \approx 0$

Classification (Binary Classification)

$$y = \{0, 1\} \quad 0 : \text{negative}$$

1 : positive

$$x = \mathbb{R}^n \rightarrow \text{not always easy!!!}$$

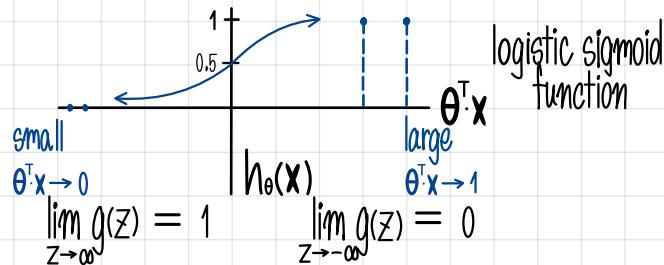
\uparrow Documents : term frequencies
fixed length List of credit card transaction

vector need to aggregate variable
length input into a "summary"

Logistic Regression

$$X_1 \rightarrow = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

but! range must be limited to [0...1]
value can be threshold
($> 0.5 \rightarrow 1, < 0.5 \rightarrow 0$)



$$h_\theta(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$$

still have parameter
 $\theta_0, \theta_1, \theta_2, \dots, \theta_n$
how to learn them?

$$g(0) = \frac{1}{1+e^0} = \frac{1}{2}$$

$$g(1000) = \frac{1}{1+e^{-1000}} = 1 \quad g(-1000) = \frac{1}{1+e^{+1000}} = 0$$

can we use squared error?

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Yes, but negative log likelihood works better

$$L(\theta) = P(y | X; \theta) \rightarrow \text{independence of training set items}$$

$$= \prod_{i=1}^m P(y^{(i)} | X^{(i)}; \theta)$$

so what is $P(y | X; \theta)$

Assume

$$P(y | X; \theta) = \begin{cases} h_\theta(x), & y = 1 \\ 1 - h_\theta(x), & y = 0 \end{cases}$$

$$\text{Example, } x^{(i)} = \begin{bmatrix} 1.2 \\ 3.1 \\ 0.5 \end{bmatrix} \quad y^{(i)} = 0$$

$$h_\theta(x^{(i)}) = \frac{1}{1+e^{[\theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \theta_3 x_3^{(i)}]}}$$

Suppose $h_\theta(x^{(i)}) = 0.9$ bad pred $P(y | X; \theta) = 1 - h_\theta(x)$
Prediction is it good or bad?

Bad! but the model says $P(y^{(i)} = 0 | X^{(i)}; \theta) = 0.1$
According to the model, the probability of the correct prediction is small

Maximum likelihood says "let's make the model's predicted probability for the correct label as high as possible"

We will try to maximize

$$P(y | X; \theta) = \prod_{i=1}^m P(y^{(i)} | X^{(i)}; \theta)$$

which is maximized by θ maximizing

$$\log P(y | X; \theta) = \log \prod_{i=1}^m P(y^{(i)} | X^{(i)}; \theta)$$

$$\log L(\theta) = \sum_{i=1}^m \log P(y^{(i)} | X^{(i)}; \theta)$$

How to maximize this when we have a case statement ??

In general, if we have

$$A = \begin{cases} B & \text{if } C = 0 \\ D & \text{if } C = 1 \end{cases} \therefore A = B^{(1-C)} D^C$$

$$\text{Therefore, } P(y^{(i)} | X^{(i)}; \theta) = h_\theta(x^{(i)})^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{(1-y^{(i)})}$$

Therefore, θ maximizing $L(\theta)$ maximizes

$$\begin{aligned} l(\theta) &= \log L(\theta) = \sum_{i=1}^m \log P(y^{(i)} | X^{(i)}; \theta) \\ &= \sum_{i=1}^m [\log h_\theta(x^{(i)})^{y^{(i)}} + \log (1 - h_\theta(x^{(i)}))^{(1-y^{(i)})}] \\ &= \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1 - h_\theta(x^{(i)})) \end{aligned}$$

$$\text{Let } \hat{y}^{(i)} = h_\theta(x^{(i)}) \quad \log x^y = y \log x$$

$$= \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1 - \hat{y}^{(i)})$$

which is maximized by θ minimizing

$$J(\theta) = - \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1 - \hat{y}^{(i)})$$

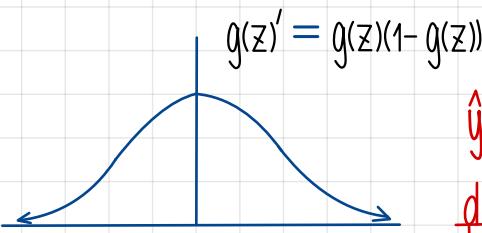
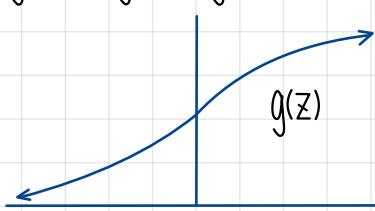
Binary Cross Entropy Cost Function

$$\hat{y}^{(i)} = \frac{1}{1+e^{-\theta^T x^{(i)}}}$$

what is the derivative of sigmoid function?

$$\frac{df}{dz} = \frac{gf' - fg'}{g^2} \quad \frac{d}{dz}(1+e^{-z}) = -e^{-z}$$

$$\begin{aligned} g(z)' &= \frac{d}{dz} \frac{-1}{(1+e^{-z})} \\ &= \frac{(1+e^{-z}) \cdot 0 - 1 \cdot (-e^{-z})}{(1+e^{-z})^2} \\ &= \frac{e^{-z}}{(1+e^{-z})^2} \\ &= \frac{1}{(1+e^{-z})} \frac{1+e^{-z}-1}{(1+e^{-z})} \\ g(z)' &= g(z)(1-g(z)) \end{aligned}$$



$$\begin{aligned} \nabla_J(\theta) &= \frac{\partial}{\partial \theta_j} J(\theta) \\ &= \frac{\partial}{\partial \theta_j} \left[- \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right] \\ &= \frac{\partial}{\partial \theta_j} \left[y \log \hat{y} + (1-y) \log (1-\hat{y}) \right] \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \frac{1}{1+e^{-\theta_j x^{(i)}}} &= \frac{(1+e^{-\theta_j x^{(i)}})' 1 - 1'(1+e^{-\theta_j x^{(i)}})}{(1+e^{-\theta_j x^{(i)}})^2} = \frac{e^{-\theta_j x^{(i)}}}{(1+e^{-\theta_j x^{(i)}})^2} (x_j) \\ &= \frac{1+e^{-\theta_j x^{(i)}} - 1(x_j)}{(1+e^{-\theta_j x^{(i)}})^2} = \frac{1}{(1+e^{-\theta_j x^{(i)}})} \left(\frac{(1+e^{-\theta_j x^{(i)}})}{(1+e^{-\theta_j x^{(i)}})} - \frac{1}{(1+e^{-\theta_j x^{(i)}})} \right) (x_j) \\ &= \frac{(1+e^{-\theta_j x^{(i)}})(1+e^{\theta_j x^{(i)}})}{(1+e^{-\theta_j x^{(i)}})^2} (x_j) \\ &= \frac{y}{\hat{y}} (1-\hat{y}) (x_j) + (1-y) \frac{1}{1-\hat{y}} (-\hat{y}(1-\hat{y})) (x_j) \\ &= y(1-\hat{y})(x_j) + (1-y)(-\hat{y})(x_j) \\ &= (y - y\hat{y} - \hat{y} + y\hat{y})(x_j) = (y - \hat{y})(x_j) \end{aligned}$$

$$\nabla_J(\theta) = (y - h_\theta(x)) x_j$$

Stochastic Gradient Descent

$$\nabla_l(\theta) = \begin{bmatrix} \frac{\partial l}{\partial \theta_0} \\ \frac{\partial l}{\partial \theta_1} \\ \vdots \\ \frac{\partial l}{\partial \theta_n} \end{bmatrix} = (y - h_\theta(x)) x$$

Ascent

Batch Gradient Descent

$$\nabla_l(\theta) = \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x^{(i)}$$

$$\nabla_l(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

Gradient Descent Rule (Batch)

$$\theta^{i-1} \leftarrow \theta^i - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

$\frac{\partial l(\theta)}{\partial \theta_j}$

Gradient Descent Rule (Stochastic)

$$\theta^{i-1} \leftarrow \theta^i - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

(update with one pattern)

Classification

1. Dataset $x \in \mathbb{R}^n \quad y \in \{0, 1\}$

2. Hypothesis function $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$

3. Cost function $J(\theta) = - \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})$

4. Optimization algorithm

$\theta^{(0)} : \text{initialize arbitrarily}$

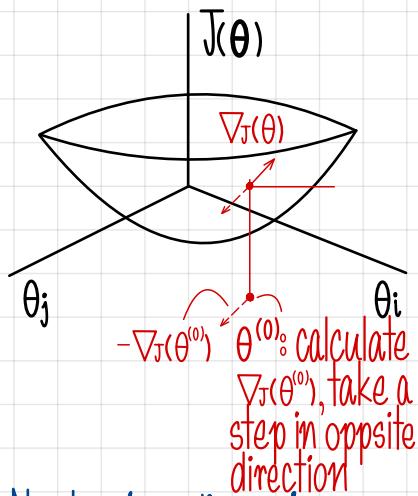
while not converged

$$\theta^{n-1} \leftarrow \theta^n - \alpha \nabla_J(\theta^{(j)})$$

J is convex \rightarrow Gradient Descent always converges if α small enough

Newton's method

Gradient Descent is the first order method



Issues :

(1) learning rate difficult to find, especially when no normalization

(2) How many iteration

the second order method can use the shape of $J(\theta)$ at current θ to decide what step

$$\text{Calculate Gradient: } \theta^{(1)} \leftarrow \theta^{(0)} - \frac{J'(\theta^{(0)})}{J''(\theta^{(0)})}$$

$$J(\theta) = J(\theta^{(0)}) = 2\theta - 6$$

$$J'(6) = 6 - 10g = 5$$

$$J''(\theta) = 2$$

$$\theta^{(1)} = 4\frac{1}{3}$$

$$J'(4\frac{1}{3}) = 9.48$$

$$J''(4\frac{1}{3}) = 21.3$$

$$\theta^{(2)} = 4\frac{1}{3} - \frac{9.48}{21.3}$$

$$\theta^{(10)} = 3.1$$

$$\theta^{(11)} = 3.067$$

$$= 3.06$$

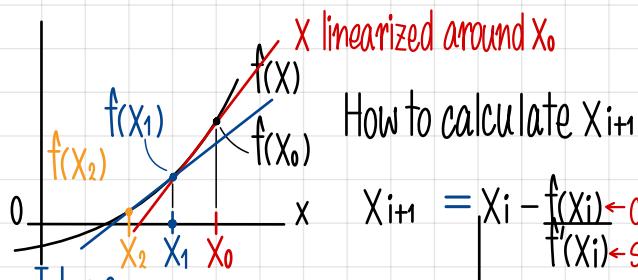
So Newton's method can be used to find the minimum of cost function.

Quickly if the function is simple & smooth

Newton in 1 dimension

zero-finding algorithm

given some complex function,
we want to find where it crosses 0



Idea:

(1) select x_0 arbitrarily

(2) let $i = 0$

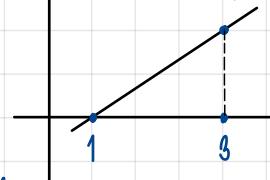
(3) calculate $f(x_i)$

(4) linearized $f(x)$ around x_i

(5) find x_{i+1} as zero of linearized $f(x)$

(6) let $i \leftarrow i + 1$

(7) go to step (3)



$$y = 1/2x - 1/2$$

want x such that y is 0

$$x = 3, y = 1$$

adjust x by height

divided by the slope

$$x_0 = 3, f(x_0) = 1$$

$$x_1 = 3, f(x_1) = 1$$

vector version of the second derivative is the Hessian.

$$\nabla^2 J(\theta) = \begin{bmatrix} \frac{\partial^2 J}{\partial \theta_0^2} & \frac{\partial^2 J}{\partial \theta_0 \partial \theta_1} & \cdots & \frac{\partial^2 J}{\partial \theta_0 \partial \theta_n} \\ \vdots & \ddots & \ddots & \ddots \\ \frac{\partial^2 J}{\partial \theta_n \partial \theta_0} & \frac{\partial^2 J}{\partial \theta_n \partial \theta_1} & \cdots & \frac{\partial^2 J}{\partial \theta_n \partial \theta_n} \end{bmatrix} \quad H_J(\theta) = \begin{bmatrix} \frac{\partial^2 J}{\partial \theta_0^2} & \frac{\partial^2 J}{\partial \theta_0 \partial \theta_1} & \cdots & \frac{\partial^2 J}{\partial \theta_0 \partial \theta_n} \\ \vdots & \ddots & \ddots & \ddots \\ \frac{\partial^2 J}{\partial \theta_n \partial \theta_0} & \frac{\partial^2 J}{\partial \theta_n \partial \theta_1} & \cdots & \frac{\partial^2 J}{\partial \theta_n \partial \theta_n} \end{bmatrix}$$

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - (H_J(\theta))^{-1} \nabla J(\theta)$$

$$\text{what } \frac{\partial^2 J}{\partial \theta_j \partial \theta_k}; \frac{\partial J}{\partial \theta_j} = (y - h_\theta(x)) x_j$$

$$\frac{\partial}{\partial \theta_k} (y - \hat{y}) x_j = x_j \frac{\partial}{\partial \theta_k} (y - \hat{y})$$

$$= x_j \frac{\partial}{\partial \theta_k} (\hat{y})$$

$$\hat{y} = \frac{1}{1 + e^{-\theta^T x}}$$

$$= -x_j \hat{y} (1 - \hat{y}) \frac{\partial}{\partial \theta_k} (\theta^T x)$$

$$= -x_j (\hat{y} - \hat{y}^2) x_k$$

$$= -(\hat{y} - \hat{y}^2) x_j x_k$$

$$H_J(\theta) = -(\hat{y} - \hat{y}^2) X \cdot X^T$$

Newton's method find a zero.

In optimization, we don't want a zero of the Cost function. We want the minimum!

We can use Newton's method to find a zero of the Gradient of J
find zero of $\rightarrow \nabla J(\theta)$
this function

$$\begin{aligned} J(\theta) &= (\theta - 3)^2 + 2 \\ &= \theta^2 - 6\theta + 9 + 2 \\ &= \theta^2 - 6\theta + 11 \end{aligned}$$

start at $\theta^{(0)} = 6$

use Newton to find $\theta^{(1)}$

$$J(\theta) = (\theta - 3)^4 + 2 \quad (\text{min is } \theta = 3)$$

$$J'(\theta) = 4(\theta - 3)^3$$

$$J''(\theta) = 12(\theta - 3)^2$$

$$\nabla J(\theta^{(0)}) = 10g$$

$$\nabla^2 J(\theta^{(0)}) = 10g$$

What can go wrong ??

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - H_J^{-1}(\theta^{(i)}) \cdot \nabla_J(\theta^{(i)})$$

rank of $H_J(\theta)$

$\hookrightarrow n+1$: full rank, invertible
 $< n+1$: rank deficient, uninvertible

$$\text{rank}(XX^T) = 1$$

$$\text{rank}(X^{(1)}X^{(1)T} + X^{(2)}X^{(2)T}) = 2$$

as long as $X^{(1)}$ and $X^{(2)}$ are linearly independent

If $m \leq n$, maximum rank of H is n .

Newton's method requires at least $n+1$

linearly independent training patterns

Generalized Linear Models (GLMs)

up to now :

Linear Regression

$\theta^T x$ is a predictor of the mean of a gaussian distribution

$$y \sim N(\theta^T x, \sigma^2)$$

Logistic Regression

$\frac{1}{1+e^{-\theta^T x}}$ is a predictor of the probability that x comes from class 1, $(1 - \frac{1}{1+e^{-\theta^T x}})$ is a predictor of the probability that x comes from class 0

$$y \sim \text{Bernoulli}\left(\frac{1}{1+e^{-\theta^T x}}\right)$$

Example Bernoulli: fair coin toss chance of the mean is Bernoulli(1/2)

Both gaussian and bernoulli distributions are members of the exponential family of distributions, any members of the family can be written in canonical form

$$P(y; \eta) = b(y)e^{(\eta^T T(y) - a(\eta))}$$

η is the natural parameter of the distributions

$T(y)$ is the sufficient ($T(y) = y$ for us)

$a(\eta)$ is the log partition function (Normalizer)

$b(y)$ is an arbitrary scalar function

Bernoulli(ϕ) is a member of the exponential family

$$P(y=1; \phi) = \phi \quad P(y=0; \phi) = 1 - \phi$$

$$P(y; \phi) = \phi^y (1-\phi)^{1-y}$$

$$\text{substitute } z = e^{\log z} = e^{\log(\phi^y (1-\phi)^{1-y})} = e^{\log \phi^y + \log(1-\phi)^{1-y}}$$

$$\text{substitute } z = e^{\log z} = e^{y \log \phi + (1-y) \log(1-\phi)} = e^{y \log \phi + \log(1-\phi) - y \log(1-\phi)}$$

$$= e^{(\log \phi - \log(1-\phi))y + \log(1-\phi)}$$

$$b(y)e^{(\eta^T T(y) - a(\eta))} = e^{(\log \frac{\phi}{1-\phi})y + \log(1-\phi)}$$

$$\eta = \log \frac{\phi}{1-\phi} \quad T(y) = y \quad b(y) = ? \quad a(\eta) = ? \quad (\text{scalar } \eta^T = \eta)$$

$$\eta = \log \frac{\phi}{1-\phi}$$

$$\eta = \log \phi - \log(1-\phi)$$

$$\log \phi = \eta + \log(1-\phi)$$

$$\phi = e^{\eta + \log(1-\phi)}$$

$$= e^{\eta} e^{\log(1-\phi)}$$

$$= e^{\eta} (1-\phi)$$

$$\phi = e^{\eta} - e^{\eta} \phi$$

$$\therefore \phi = e^{\eta} / (1+e^{\eta})$$

this prove that Bernoulli distribution with parameter ϕ is in exponential family

Gaussian distribution with $\sigma^2 = 1$ is a member of the exponential family

$$\eta = \mu, T(y) = y, a(\eta) = \frac{\mu^2}{2} = \frac{\eta^2}{2}$$

$$b(y) = \frac{1}{\sqrt{2\pi\sigma^2}}$$

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2} + \mu y - \frac{\mu^2}{2}}$$

What is this GLMs anyway?

Assumes 3 things

1. $P(y|x; \phi)$ = exponential family(η)

2. Given X , we want to predict an expected value of $T(y)$ (or y in most cases) Given X

3. η is linear in x , i.e.,

$$\eta_i = \theta_i^T x$$

expected \rightarrow Logistic Regression

Logistic Regression

$$E(y|x) = 0 \cdot P(y=0|x) + 1 \cdot P(y=1|x)$$

$$\phi = \frac{e^\eta}{1+e^\eta} = \frac{1}{1+e^{-\eta}}$$

$$= P(y=1|x)$$

[ϕ]

GLMs principle says that if we want to predict ϕ for a Bernoulli distributions as $P(y=1|x)$, then we should use $P(y=1|x) = 1/(1+e^{-\theta x})$

For Linear Regression, predict

$$E(y) = \mu = \theta^T x$$

why? Because $\eta = \mu$ give us the exponential family representation of the normal distribution.

For Logistic Regression, predict

$$E(y) = P(y=1|x) = 1/(1+e^{-\theta x})$$

why? Because $\eta = \log \phi / (1-\phi)$ give us the exponential family representation of the bernoulli distribution.

GLM procedure :

write target y as result of a probabilistic experiment.

$$y \sim P(\phi)$$

convert $P(\phi)$ to a member of the exponential family with parameter η .

Plug $\eta_i = \theta^T x$ into the expression for ϕ

This is give you a Hypothesis function $h_{\theta}: x \rightarrow y$

Cost function is always negative log Likelihood

$$J(\theta) = l(\theta) = P(x^{(i)}|y^{(i)}, \theta)$$

Take $\nabla J(\theta) = 0$, use analytical solution or Gradient Descent or Newton's method or...?

New Distribution : Multinomial

Given x , we would like to predict the class of x , where the class is exactly one of K categories, we assume

$$y \sim \text{Multinomial}(\phi_1, \phi_2, \dots, \phi_K)$$

Multinomial Regression or

Multinomial Logistic Regression

$$h_{\theta}(x) = \begin{bmatrix} \frac{e^{\theta_1^T x}}{\sum_{j=1}^K e^{\theta_j^T x}} \\ \vdots \\ \frac{e^{\theta_K^T x}}{\sum_{j=1}^K e^{\theta_j^T x}} \end{bmatrix}$$

Soft max Function

Cross Entropy Cost Function

Generative Learning algorithm

Sofar, **Discriminative** models that learn $P(y|x)$ directly

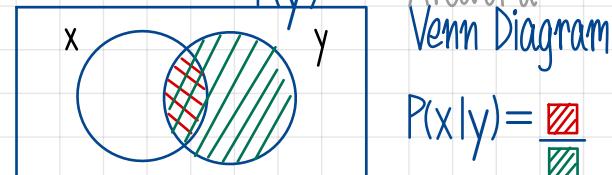
Logistic Regression $P(y=1|x; \theta) = 1/(1+e^{-\theta x})$
never consider $P(x|y)$ or $P(x)$ or $P(y)$

Alternative (**Generative**) approach
use that fact (**Baye's Rule**)

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

(why? Definition of condition)

$$\text{tellus } P(x|y) = \frac{P(x \wedge y)}{P(y)} = \frac{\text{Area of } x \cap y}{\text{Area of } y}$$



$$P(x|y) = \frac{\text{Red Area}}{\text{Green Area}}$$

$$P(x|y) = \frac{P(x \wedge y)}{P(y)} \rightarrow P(x \wedge y) = P(x|y)P(y)$$

$$P(y|x) = \frac{P(y \wedge x)}{P(x)} \rightarrow P(y \wedge x) = P(y|x)P(x)$$

$$P(x|y)P(y) = P(y|x)P(x)$$

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

Baye's Rule

If we want, given a feature vector x , to estimate a probabilistic distribution over y , we can calculate $P(y|x)$ directly (Logistic Regression) or as $\frac{P(x|y)P(y)}{P(x)}$

Discriminative Generative

Generative approach supposes that we actually model the process by which the dataset was generated

for $i = 1 \text{ to } n$

$$y^{(i)} \sim \{0, 1\} \text{ according to prior distribution } P(y)$$

Given x , what is the maximum a posterior estimate of y

$$y^* = \operatorname{argmax}_y P(y|x)$$

$$= \operatorname{argmax}_y \frac{P(x|y)P(y)}{P(x)}$$

$$= \operatorname{argmax}_y P(x|y)P(y)$$

$P(y|x)$ posterior

$P(x|y)$ class-conditional probability

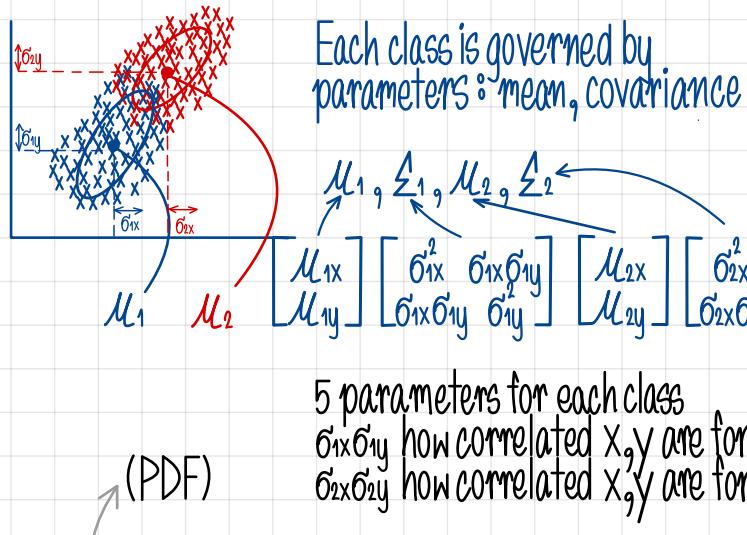
$P(y)$ prior

$P(x)$ data probability (useless constant)

Example $P(x|y)$ is multivariate gaussian

$Y = \{\text{male, female}\}$

$X = \mathbb{R}^2$ two features: Person's height and weight



Covariance matrix specifies elliptical ISO contours of the probabilistic distribution function. Major axis of the elliptical contours. For one class is the eigenvector of the covariance matrix corresponding to its largest eigenvalue.

Given all this,

$$P(x|y=j) = \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} e^{-\frac{1}{2}(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}$$

$$= N(\mu_j, \Sigma_j)$$

Suppose height = 165, weight = 65

What is y^* ? $y^* = \underset{y}{\operatorname{argmax}} P(x|y) P(y)$

Suppose I calculate

$$P(x|y=1) = 0.12 \quad P(x|y=2) = 0.10$$

Suppose in general population $P(y) = \phi^y (1-\phi)^{1-y}$

$$\begin{aligned} \text{priors } P(y=1) &= 0.51 \quad P(y=2) = 0.49 \\ P(y=1|x) &\sim 0.12 \cdot 0.51 \\ P(y=2|x) &\sim 0.10 \cdot 0.49 \end{aligned}$$

Suppose the sample is among students DSAI program at AIT

$$\begin{aligned} \text{female } P(y=1) &= 0.4 \quad P(y=2) = 0.6 \\ P(y=1|x) &\sim 0.12 \cdot 0.4 = 0.048 \\ P(y=2|x) &\sim 0.10 \cdot 0.6 = 0.06 \end{aligned}$$

$$P(x|y=1) = 0.12 \quad P(x|y=2) = 0.10$$

$$P(y=1) = 0.4 \quad P(y=2) = 0.6 \quad P(x) ???$$

$$P(y=1|x) + P(y=2|x) = 1$$

$$\frac{P(x|y=1)P(y=1)}{P(x)} + \frac{P(x|y=2)P(y=2)}{P(x)} = 1$$

$$P(x) = P(x|y=1)P(y=1) + P(x|y=2)P(y=2)$$

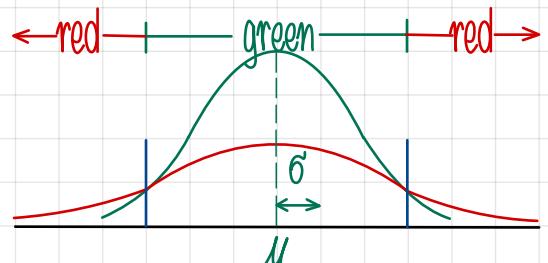
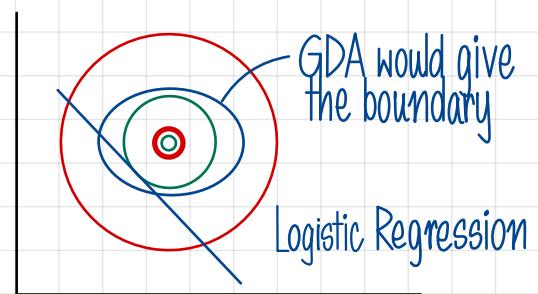
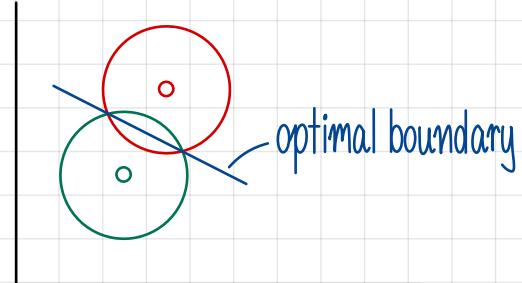
$$P(x) = \sum_{y \in Y} P(x|y=y_i)P(y=y_i)$$

$$P(x=[165, 65]) = 0.12 \cdot 0.4 + 0.10 \cdot 0.6 = 0.108$$

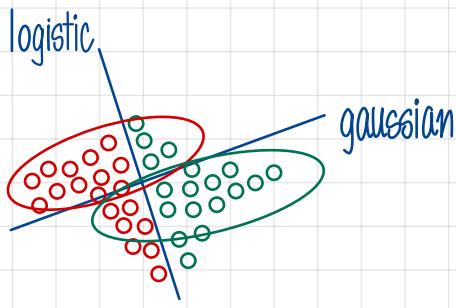
Generative Model with Multivariate Gaussian

GDA: Gaussian Discriminant Analysis

If $P(y=1) = P(y=2)$, $\Sigma_1 = \Sigma_2$, GDA will give the same boundary as Logistic Regression



If you know that the data for each class are approximately gaussian and covariance are not equal, gaussian generative model will outperform the logistic regression model.



$$P(y|x) \propto P(x_1|y)P(x_2|x_1,y)P(x_3|x_1,x_2,y)P(x_4|x_1,x_2,x_3,y)P(y)$$

$4 \times 3 \times 5 \times 2$
360 entries in the table!

Naive Bayes Assumption:

$$P(x_j|x_1, \dots, x_{i-1}, y) \approx P(x_j|y)$$

$$P(\text{Colors}=\text{red} | \text{Size}=\text{medium}, \text{Outcome}=\text{Buy}) \approx$$

$$P(\text{Colors}=\text{red} | \text{Outcome}=\text{Buy})$$

$$P(y|x) \propto P(x_1|y)P(x_2|y)P(x_3|y)P(x_4|y)P(y)$$

approx

Complicated Joint Distribution is approximated by a product of simple conditional probabilities

Advantage: fewer parameters

Parameters:

$$x_i = \{S, M, L\}$$

$$P(x_i|y)$$

x_i	B	$\neg B$
S	0.3	0.3
M	0.3	0.5
L	0.4	0.2
total	1.0	1.0

$$y_i = \{B, \neg B\}$$

$$\text{Size}$$

$$\text{Colors}$$

$$\text{Gender}$$

$$\text{Age}$$

$$\text{Outcome}$$

6 parameters

8 parameters

6 parameters

10 parameters

2 parameters

$$\begin{matrix} P(X_2|y) \\ P(X_3|y) \\ P(X_4|y) \\ P(y) \end{matrix}$$

32 parameters total!

vs 362 for exact calculation!

$$180 \times 2 + 2 \quad P(x|y)P(y)$$

Generally speaking, A simple model will generalize better (less chance to overfit). Naive Bayes approximation is a **Regularizer**. Naive Bayes usually much better than full bayesian Classification.

How many parameters in Gaussian model?
mean, covariance

$$X = \mathbb{R}^2$$

mean k vectors of length N, $k \cdot N$ parameters

$$Y = \{1, \dots, k\}$$

$$n, k$$

covariance

$$\begin{bmatrix} 1 & 1 & 1 & \dots \end{bmatrix}$$

$$\sum_{i=1}^n = \frac{n(n+1)}{2}$$

$$k \frac{n(n+1)}{2}$$

2 classes, 10 parameters

$$\frac{2 \cdot 10}{2 \cdot 10 \cdot 11}$$

(dominated by kN^2)

$$20 + 110 = 130 \quad (40 \text{ for Diagonal covariance matrices})$$

$$j \quad i \quad [\quad]$$

covariance between x_i, x_j

$$P(b|a) = P(a, b)$$

$$P(b)$$

$$P(b|a)P(b) = P(a, b)$$

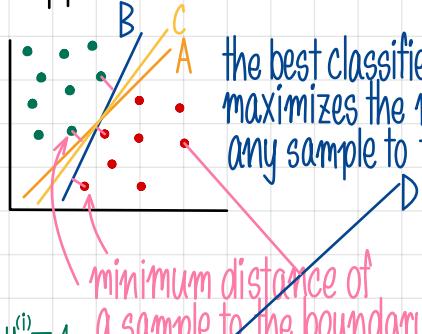
$$P(a)$$

$$P(a)$$

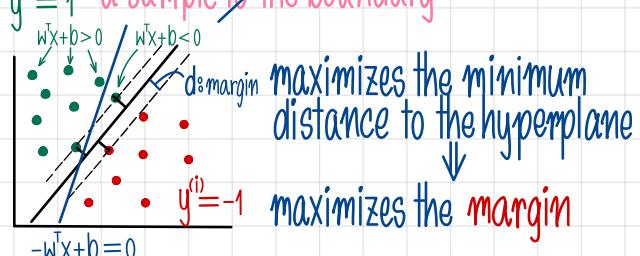
$$P(a, b)$$

$$P(a)$$

Support Vector Machines (SVM)



the best classifier is the boundary that maximizes the minimum distance of any sample to the boundary



$$\text{Line in 2D: } ax_1 + bx_2 + c = 0$$

$$\text{Line in 3D: } ax_1 + bx_2 + cx_3 + d = 0$$

weights bias
w₁ w₂ w₃ b

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \quad b$$

parameters to estimate: w, b

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \text{ "margin"}$$

parameters w, b define the separating hyperplane

$$\text{Assume } Y = \{-1, 1\}$$

$$X = \mathbb{R}^2$$

$$h_{w,b}(x) = g(w^T x + b) \quad \text{logistic Regression}$$

$$g(z) = \begin{cases} 1 & z \geq 0 \\ -1 & \text{otherwise} \end{cases} \quad \text{augmented with 1}$$



Margin

$$\text{Functional margin: } y^{(i)} = 1 \quad w^T x + b > 0$$

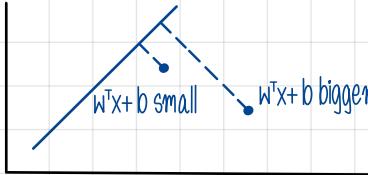
$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b)$$

+1 ≥ 0 for A+1 prediction
-1 < 0 for A-1 prediction

If $\hat{\gamma}^{(i)} > 0$ for $i \in 1 \dots m$ $h(x^{(i)}) = y^{(i)}$
All training examples are correctly classified.

If $\hat{\gamma}^{(i)} < 0$ for any i, then samples i is incorrectly classified

$\hat{\gamma}^{(i)}$ increases as $x^{(i)}$ is moved further from the boundary



We said our goal was to maximize the margin
First attempt (incorrect!!)

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \min_i \hat{\gamma}^{(i)}$$

Problem: $\|w\|$ is not bounded range of

If w, b give $\min_i \hat{\gamma}^{(i)} = \varepsilon > 0$,

then sw, sb give $\min_i \hat{\gamma}^{(i)} = s\varepsilon > 0$ where $s > 1$

Instead of w, b we would like to normalize:

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \min_i \left[\frac{y^{(i)}}{\|w\|} \left(\frac{w^T x^{(i)}}{\|w\|} + \frac{b}{\|w\|} \right) \right]$$

Geometric margin: of hyperplane (w, b)

$$\gamma^{(i)} = y^{(i)} \left[\left(\frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right]$$

Now we can revise the objection to

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \min_i \gamma^{(i)}$$

$$\text{Let } \gamma = \min_i \gamma^{(i)} \rightarrow w^*, b^* = \underset{w, b}{\operatorname{argmax}} \gamma$$

Goal (Revised)

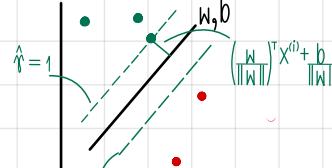
$$\max_{\gamma, w, b} \gamma$$

$$\text{subject to } y^{(i)}(w^T x^{(i)} + b) \geq \gamma, i=1 \dots m \quad \|w\| = 1$$

$$\max_{\gamma, w, b} \gamma \quad \left\{ \begin{array}{l} \text{non-convex} \\ \gamma = \frac{\hat{\gamma}}{\|w\|} \end{array} \right.$$

$$\text{subject to } y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, i=1 \dots m$$

→ instead of requiring $\|w\| = 1$ force $\|w\| = \text{whatever give us } \hat{\gamma} = 1$



$$\hat{\gamma} = \min_i \hat{\gamma}^{(i)}$$

$$\hat{\gamma} = y^{(i)}(w^T x^{(i)} + b)$$

If $\gamma = 1.0$, for example

If we let $w \leftarrow \frac{w}{\|w\| \cdot 1.0}$ $b \leftarrow \frac{b}{\|w\| \cdot 1.0}$

$$\hat{\gamma} = 1 \\ \min_{w, b} \frac{1}{2} \|w\|^2$$

subject to $y^{(i)}(w^T x^{(i)} + b) \geq 1, i=1 \dots m$

$$\frac{1}{2} \|w\|^2 = \frac{1}{2}(w_1^2 + w_2^2 + w_3^2 + \dots + w_n^2)$$

Objection Function Quadratic in parameters
constraints that are linear, QP
Quadratic Programming
QP given

$$c \in \mathbb{R}^n \quad Q \in \mathbb{R}^{n \times n}$$

$$A \in \mathbb{R}^{m \times n} \quad b \in \mathbb{R}^m$$

Find $x^* \in \mathbb{R}^n = \arg \min_x \frac{1}{2} x^T Q x + c^T x$

subject to $Ax \leq b$

QP

$$x$$

SVM

$$\longleftrightarrow$$

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ b \end{bmatrix}$$

$$\min \frac{1}{2} x^T Q x + c^T x \longleftrightarrow \min \frac{1}{2} (w_1^2 + \dots + w_n^2)$$

$$Q = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix} \quad c = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$$

$$y^{(i)}(w^T x^{(i)} + b) \geq 1$$

$$\begin{bmatrix} y^{(1)} x_1^{(1)} & y^{(1)} x_2^{(1)} & \dots & y^{(1)} x_n^{(1)} & y^{(1)} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ b \end{bmatrix} \geq [1]$$

$$\begin{bmatrix} -y^{(1)} x_1^{(1)} & -y^{(1)} x_2^{(1)} & \dots & -y^{(1)} x_n^{(1)} & -y^{(1)} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ b \end{bmatrix} \leq [-1]$$

$$A = \begin{bmatrix} -y^{(1)} x^{(1)\top} & -y^{(1)} \\ \vdots & \vdots \\ -y^{(m)} x^{(m)\top} & -y^{(m)} \end{bmatrix} \quad b = \begin{bmatrix} -1 \\ \vdots \\ -1 \end{bmatrix}$$

To find w, b maximizes the margin, let

$$Q = \begin{bmatrix} I_{n \times n} & 0 \\ 0 & 0 \end{bmatrix}$$

$$c = 0_{n \times 1} \quad \begin{bmatrix} x_1^{(1)} & \dots & x_n^{(1)} \\ \vdots & & \vdots \\ x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

$$b = -1_{m \times 1}$$

$$A = -y^* x \leftarrow$$

pip install cvxopt \leftarrow "convex optimization"

\odot We can maximizes the margin!

\odot There must be a linear hyperplane "shattering" the training set

\odot Not for big m or big n !
 1920×1080 ?

$$x^T Q$$

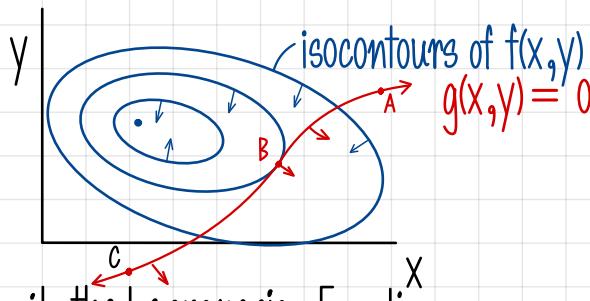
$$\begin{bmatrix} w_1 & w_2 & \dots & w_n & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ 0 \end{bmatrix} \quad w_1^2 + \dots + w_n^2 + 0$$

$$Ax \leq b$$

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & & & \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \leq \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

Lagrangian Optimization

Example: maximize $f(x, y)$ subject to $g(x, y) = 0$
constrained maximization



We write the Lagrangian Function containing a Lagrange multiplier λ

$$L(x, y, \lambda) = f(x, y) - \lambda g(x, y)$$

Theorem: If (x^*, y^*) is a maximum of $f(x, y)$

satisfying $g(x^*, y^*) = 0$, then there exist λ^* such that (x^*, y^*, λ^*) is a stationary point of $L(x, y, \lambda)$

That means that (x^*, y^*, λ^*) satisfies

$$\frac{\partial L}{\partial x} = 0, \quad \frac{\partial L}{\partial y} = 0, \quad \frac{\partial L}{\partial \lambda} = 0$$

To find max of $f(x, y)$ such that $g(x, y) = 0$

(1) write the Lagrangian

(2) set derivative w.r.t x, y, λ to 0

(3) solve !!

(4) check whether the solution is a max or not

$$\begin{aligned} y &= (x-3)^2 + 2 \\ y' &= 2(x-3) \\ 0 &= 2x-6 \\ x^* &= 3 \\ y'' &= 2 \leftarrow \text{positive 2nd derivative indicates } x^* \text{ is a minimum} \end{aligned}$$

For a multivariate Function,

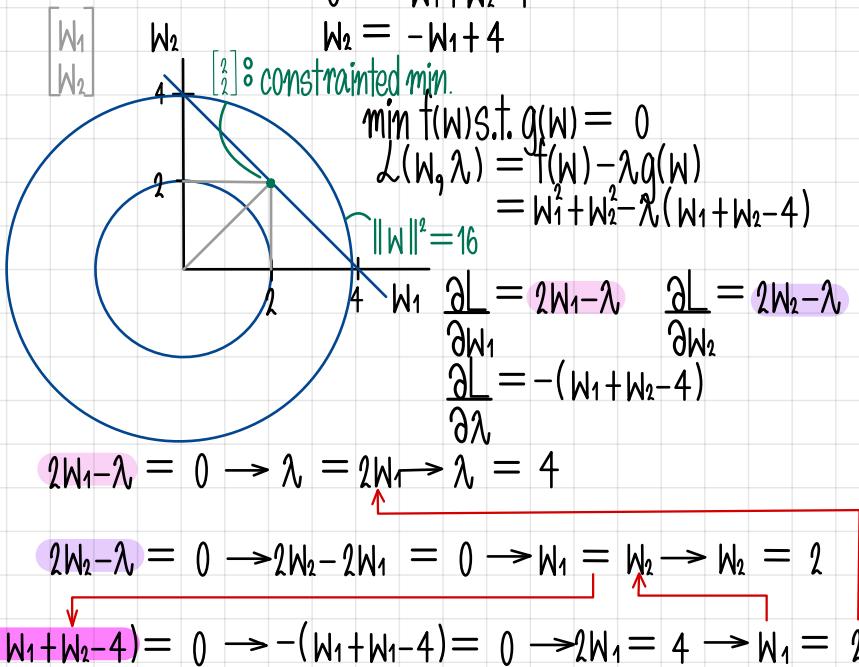
Hf: Hessian of f . Unconstrained optimization positive eigenvalues for minima,

negative eigenvalues for maxima

H_L: Constrained optimization all stationary points are "saddle" points

Specific example:

$$\begin{aligned} f(W) &= \|W\|^2, g(W) = W_1 + W_2 - 4 \\ 0 &= W_1 + W_2 - 4 \\ W_2 &= -W_1 + 4 \end{aligned}$$



$$W_1 = 2 \quad W_2 = 2 \quad \lambda = 4$$

$$\begin{array}{ll} \min_w f(w) & \\ \text{s.t. } g_i(w) \leq 0, i \in 1 \dots k & \text{If } w \text{ satisfies constraint} \\ h_i(w) = 0, i \in 1 \dots l & h_i(w) = 0, g_i(w) \leq 0 \end{array}$$

The Generalized Lagrangian

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

Primal Optimization

$$\begin{aligned} \theta_p(w) &= \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta) \\ &= \max_{\alpha, \beta, \alpha_i \geq 0} f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w) \end{aligned}$$

Claim: If w violates any constraint, then $\theta_p(w) = 0$

If w satisfies constraint, then $\theta_p(w) = f(w)$

Sofar, we only max w.r.t α, β , w is fixed ($\theta_p(w) = f(w)$ or ∞)

Let w^* be the w satisfying

$$\min_w \theta_p(w) = \min_w \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta)$$

$p^* = \min_w \theta_p(w)$ is called the value of the **primal** problem

Primal is what we want to solve, but it's not always. Instead, sometimes, the **dual** problem is easier.

$$\theta_d(\alpha, \beta) = \min_w \mathcal{L}(w, \alpha, \beta)$$

Dual problem :

$$\max_{\alpha, \beta, w \geq 0} \theta_d(\alpha, \beta) = \max_{\alpha, \beta, \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta)$$

$$d^* = \max_{\alpha, \beta, \alpha_i \geq 0} \theta_d(\alpha, \beta)$$

is called the value of the **dual** problem

Under Specific Circumstances

$$p^* = d^*$$

$p^* = d^*$ Under these conditions

f is convex

The g_i 's are convex

The h_i 's are affine

The g_i 's are strictly feasible

For SVM linearly separable data
(\exists w.s.t. $g_i(w) \leq 0$ for all)

When these conditions hold, $\exists w^*, \alpha^*, \beta^*$

(1) w^* is the solution of the **primal** problem

(2) α^*, β^* are the solution of the **dual** problem

(3) $p^* = d^* = \mathcal{L}(w^*, \alpha^*, \beta^*)$

(4) w^*, α^*, β^* satisfy the **KKT** conditions

Lagrangian Duality

$$\frac{\partial \mathcal{L}(w^*, \alpha^*, \beta^*)}{\partial w_i} = 0, i \in 1 \dots n$$

$$\frac{\partial \mathcal{L}(w^*, \alpha^*, \beta^*)}{\partial \beta_i} = 0, i \in 1 \dots l$$

$$\alpha_i^* g_i(w^*) = 0, i \in 1 \dots k \quad \text{for every } i, \alpha_i^* = 0 \text{ or } g_i(w^*) = 0$$

$$g_i(w^*) \leq 0, i \in 1 \dots k$$

$$\alpha_i^* \geq 0, i \in 1 \dots k$$

Back to SVMs!

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1, i = 1 \dots m$$

inequalities, affine g_i 's

$$g_i(w) = -y^{(i)}(w^T x^{(i)} + b) + 1 \leq 0$$

We will have one α_i for each training sample

KKT tell us $\alpha_i = 0$ or $g_i(w) = 0$

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1]$$

Generalized Lagrangian for the SVM maximum margin problem

Find $\frac{\partial \mathcal{L}}{\partial w_i}$

$$\frac{\partial \mathcal{L}(w, b, \alpha)}{\partial w_i} = \frac{1}{2} \sum_{j=1}^m \alpha_j y^{(j)} X^{(j)} X^{(i)}$$

$$\nabla_w \mathcal{L}(w, b, \alpha) = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_1} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial w_n} \end{bmatrix} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m \alpha_i y^{(i)} X^{(i)}_1 \\ \vdots \\ \sum_{i=1}^m \alpha_i y^{(i)} X^{(i)}_n \end{bmatrix}$$

$$= w - \sum_{i=1}^m \alpha_i y^{(i)} X^{(i)}$$

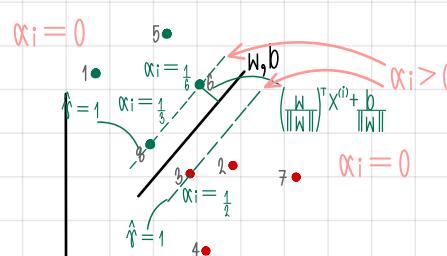
0 For $x^{(i)}$ far from boundary,
 $\alpha_i \geq 0$ For points defining boundary

$$\nabla_w \mathcal{L}(w, b, \alpha) = 0$$

$$w = \sum_{i=1}^m \alpha_i y^{(i)} X^{(i)}$$

$$\text{Find } \frac{\partial \mathcal{L}}{\partial \beta} = -\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

the sum of α_i 's for positives must be equal to the sum of α_i 's for negatives



$$S = \text{"support vectors"} = \{3, 6, 9\}$$

Plotting w^* into $\mathcal{L}(w, b, \alpha)$, we obtain

$$\begin{aligned}\mathcal{L}(w, b, \alpha) &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1] \\ &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} y^{(j)} x^{(i)} x^{(j)^T} - \sum_{i=1}^m \alpha_i [y^{(i)} (\sum_{j=1}^m \alpha_j y^{(j)} x^{(j)^T} x^{(i)}) + b) - 1] \\ &= \frac{1}{2} \left[\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right]^T \left[\sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right] - \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(j)^T} x^{(i)} - \sum_{i=1}^m \alpha_i y^{(i)} b + \sum_{i=1}^m \alpha_i \\ &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)^T} x^{(j)} - \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(j)^T} x^{(i)} - \sum_{i=1}^m \alpha_i y^{(i)} b + \sum_{i=1}^m \alpha_i \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)^T} x^{(j)} - b \sum_{i=1}^m \alpha_i y^{(i)} \\ \mathcal{L}(w, b, \alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)^T} x^{(j)}\end{aligned}$$

Generalized Lagrangian Dual problem :

$$\max_{\alpha} L(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

such that $\alpha_i \geq 0, i \in 1 \dots m$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

$$b^* = \max_{i: y^{(i)}=-1} w^* x^{(i)} + \min_{i: y^{(i)}=1} w^* x^{(i)}$$

$$b^* = \frac{1}{N_s} \sum_{i \in S} (y^{(i)} - \sum_{j \in S} \alpha_j y^{(j)} \langle x^{(i)}, x^{(j)} \rangle)$$

S : set of indices of "support vectors"
($i \text{ s.t. } \alpha_i \geq 0$)

Inference problem

given new point x what is \hat{y} ?

$$h_{w,b}(x) = \begin{cases} +1 & \text{if } w^T x^{(i)} + b \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$w^T x^{(i)} + b = \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b$$

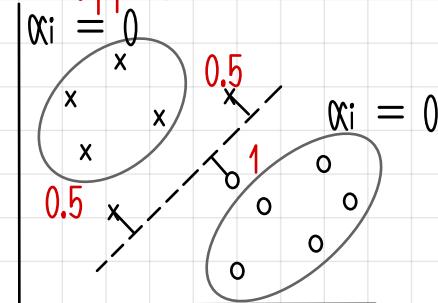
We can answer the question of which side of the boundary w, b a point x is on without ever calculating w . Only need α^*, y^*, x^* and b
(only for the S.V's)

Formulated the SVM problem
maximize the margin

Functional Geometric
Primal, Dual Lagrangian Optimization

Solution in terms of α_i 's

$\alpha_i > 0$ support vectors
 $\alpha_i = 0$ Insignificant examples



$$\max_{\alpha} L(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

such that $\alpha_i \geq 0, i \in 1 \dots m$

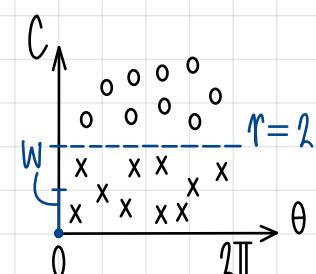
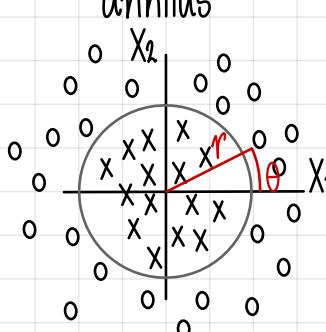
Problem

- 1 large scale
- 2 linearly separable training set
- 3 not much different from

Logistic Regression (nonlinear boundary)

Solution to Problem #3 nonlinear transformation into "Feature Space"

$$\phi: \mathbb{X} \rightarrow \mathbb{R}^2 \quad \phi = \begin{bmatrix} \text{atan2}(x_2, x_1) \\ \sqrt{x_1^2 + x_2^2} \end{bmatrix}$$



Cartesian \rightarrow Polar Idea: transform x_1, x_2 into θ, r new space where the boundary is linear

$$\max_{\alpha} L(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$

$$\begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

Possible transforms:

Polynomial

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_n \\ x_1^2 \\ \vdots \end{bmatrix}$$

Kernel function:
 $K: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$
 linear Kernel
 $K(x, x') = x_1 x'_1 + \dots + x_n x'_n = \langle x, x' \rangle$

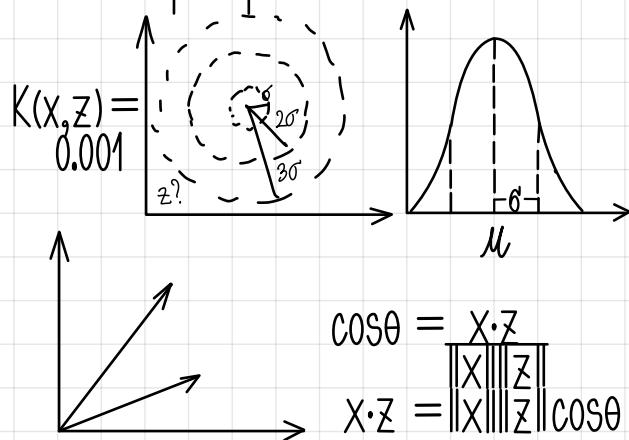
RBF (Radial Basis Function) Kernel:

$$K(x, x') = e^{-\|x - x'\|^2/2\sigma^2} \text{ or } e^{-r\|x - x'\|^2}$$

gamma

Intuition

A Kernel says how "similar" two vectors in the input space are



Linear Kernel does measure similarity in direction

Once we have our α_i 's then interference with the Kernel formulation:

$$\text{Class}(x) = \text{sgn} \left[\sum_{i \in S}^m \alpha_i y^{(i)} K(x^{(i)}, x) + b \right]$$

sum over entire training set!

Logistic Regression:

$$\text{Class}(x) = \text{sgn}(\mathbf{w}^T x + b)$$

$$K(x, z) = x_1 z_1 + \dots + x_n z_n$$

Linear: corresponds to $\phi(x) = x$

transform presentation

$$- \|x - z\|^2 / (2\sigma^2)$$

$$K(x, z) = e$$

ϕ = transforms into an infinite dimensional space!

Are there Kernel function that don't correspond to any ϕ at all? Yes

THM: Kernel function $K(x, z)$, corresponds to features space transformation ϕ IFF

The gram matrix

$$\text{is always positive semidefinite}$$

$$K = \begin{bmatrix} K(x^{(1)}, x^{(1)}) & K(x^{(1)}, x^{(2)}) & \dots \\ K(x^{(2)}, x^{(1)}) & K(x^{(2)}, x^{(2)}) & \dots \\ \vdots & \vdots & \ddots \\ K(x^{(m)}, x^{(1)}) & K(x^{(m)}, x^{(2)}) & \dots \end{bmatrix}$$

In practice

Linear

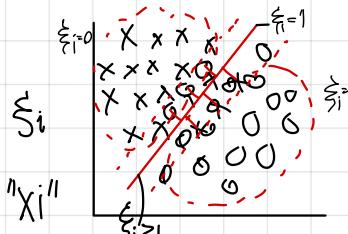
Fast (minimal Support Vector $n+1$)
 but strictly linear, worst accuracy
 least overfitting

Polynomial → Medium speed $\approx n^2$ Support Vectors
 for quadratic, etc. Medium accuracy

RBF

Medium overfitting dependent on degree
 Slow $\approx n$ Support Vectors
 Best accuracy (training set)
 Most overfitting

Solution to Problem #2 Overlap between classes



For RBF,
 Separating slack cars
 the sum has two
 Hyperparameters r, C

"SLACK" variables expressing
 How far across the boundary a point lies

Add ξ_i 's to optimization for $\epsilon > 0$

Introduce Hyperparameter

$\frac{C}{\sum_{i=1}^m \xi_i}$ strength of the "Don't overlap" heuristic compared to linearly separable
 To be minimized linearly separable data
 positive value for non margin objective

New Objective:

$$\max_{\alpha} L(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j K(x^{(i)}, x^{(j)})$$

such that $0 \leq \alpha_i \leq C$ $\sum_{i=1}^m \alpha_i y^{(i)} = 0$

Unsupervised Learning

Supervised : given (x, y) pairs

, find best $h_0 : X \rightarrow Y$

predicting target y given input x

Unsupervised : given $x, \dots ?$ **Unbalanced data**

1. Find **cluster** x belongs to

2. Masked Language Modeling (MLM)

x : "the quick fix jumped over the lazy brown dog"

x : "the $\underset{\text{quick}}{\underline{\text{fix}}}$ jumped over the lazy brown dog"

3. Dimension Reduction (Latent embedding)

$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $m \ll n$

such that $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is as accurate as possible

$$f, f' \leftarrow \underset{f, f'}{\operatorname{argmin}} \|x - f(f(x))\|$$

4. Density estimation

$f : \mathbb{R}^n \rightarrow \mathbb{R}$
 $f(x)$ predicts $p(x = x)$

K-mean clustering

input : training sets $S = \{x^{(1)}, \dots, x^{(m)}\}$

$$x^{(i)} \in X = \mathbb{R}^n$$

Goal : Group the $x^{(i)}$ into clusters

Algorithm :

1. Randomly initialize K cluster centroids

$$\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$$

2. Repeat until converge

1. for $i = 1, \dots, m$, $c^{(i)} \leftarrow \underset{j}{\operatorname{argmin}} \|x^{(i)} - \mu_j\|^2$

2. for $j = 1, \dots, K$, $\mu_j \leftarrow \frac{\sum_{i=1}^m \delta(c^{(i)} = j) x^{(i)}}{\sum_{i=1}^m \delta(c^{(i)} = j)}$

move the cluster
centroids

Gaussian mixture model



Generative model for each sample :

1. Select a component $1-K$ according to a multinomial distribution with parameters

$$\phi_1, \dots, \phi_k \text{ (priors), } \left(\sum_{i=1}^m \phi_i = 1 \right)$$

Let i be the selected component

2. sample the datum $x_j \in \mathbb{R}^n$ according to

$$x_j \sim N(\mu_i, \Sigma_i)$$

Learning Objective :

estimate priors ϕ_1, \dots, ϕ_k
means μ_1, \dots, μ_k
covariances $\Sigma_1, \dots, \Sigma_k$

Cost function : negative log likelihood

Steps :

1. write Cost function

2. take derivatives

3. try to solve the resulting system of equations

Log likelihood function

$$l(\phi, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k)$$

$$= \sum_{i=1}^m \log P(x^{(i)}, z^{(i)}; \phi, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k)$$

↳ could have been generated by any of the components in the mixture.

Let $z^{(i)}$ be the (hidden) component that generated

example i. $z^{(i)}$ could have any value 1 to K

Aside : it is always valid to write probability as a sum

of conditional probabilities

$$\begin{aligned}
 P(A) &= \sum_{b \in B} P(A|B=b)P(B=b) \\
 &\sum_{i=1}^m \log P(X^{(i)}, Z^{(i)}; \emptyset, \mu_1, \dots, \mu_k, \Sigma_1, \dots) \\
 &= \sum_{i=1}^m \log \sum_{j=1}^k P(X^{(i)}|Z^{(i)}=j; \mu_j, \Sigma_j)P(Z^{(i)}=j; \emptyset_j) \\
 &= \sum_{i=1}^m \log \sum_{j=1}^k \emptyset_j N(X^{(i)}; \mu_j, \Sigma_j) \quad X_i \text{ } \bigcirc \text{ } \mu_j \\
 &= \sum_{i=1}^m [\log P(X^{(i)}|Z^{(i)}; \mu_z, \Sigma_z) + \log P(Z^{(i)}; \emptyset)] \\
 &= \sum_{i=1}^m [\log \emptyset_z + \log N(X^{(i)}; \mu_z, \Sigma_z)]
 \end{aligned}$$

$$\frac{\partial l}{\partial \emptyset_j} = 0 \quad \sum_{i=1}^m \delta(Z^{(i)}=j) \frac{1}{\emptyset_j} = 0$$

\$\emptyset_j = \infty\$? or very large? not
 satisfactory, \$\sum_{j=1}^m \emptyset_j = 1 \rightarrow \sum_{j=1}^m \emptyset_j - 1 = 0\$

To solve the constrained Optimization

use Lagrangian

$$\begin{aligned}
 l(\emptyset, \mu_1, \dots, \mu_k, \Sigma_1, \dots) &+ \lambda \left(\sum_{j=1}^m \emptyset_j - 1 \right) \\
 \sum_{i=1}^m \frac{N(X^{(i)}; \mu_z, \Sigma_z)}{\sum_{l=1}^m \emptyset_l N(X^{(i)}; \mu_z, \Sigma_z)} + \lambda &= 0
 \end{aligned}$$

Skip sumsteps $\lambda = -m$

Also task derivatives $\frac{\partial l}{\partial \mu_j} = 0$

with some works, you get

$$\begin{aligned}
 \emptyset_j &= \frac{1}{m} \sum_{i=1}^m \delta(Z^{(i)}=j) \\
 \mu_j &= \frac{\sum_{i=1}^m \delta(Z^{(i)}=j) X^{(i)}}{\sum_{i=1}^m \delta(Z^{(i)}=j)} \\
 \Sigma_j &= \frac{\sum_{i=1}^m \delta(Z^{(i)}=j) (X^{(i)} - \mu_j)(X^{(i)} - \mu_j)^T}{\sum_{i=1}^m \delta(Z^{(i)}=j)}
 \end{aligned}$$

Just like the GDA Model

Expectation maximization (EM)

Algorithm
parameters estimation under uncertain hidden variable

Two Steps

E-Steps (Expectation) :
guess values of hidden variables ($Z^{(i)}$)

M-Steps (Maximization) :
maximum likelihood estimation (MLE)
based on the guess

EM Algorithm :

Repeat until convergence :

E-Steps :

foreach $i \in 1, \dots, m$, $j \in 1, \dots, k$,

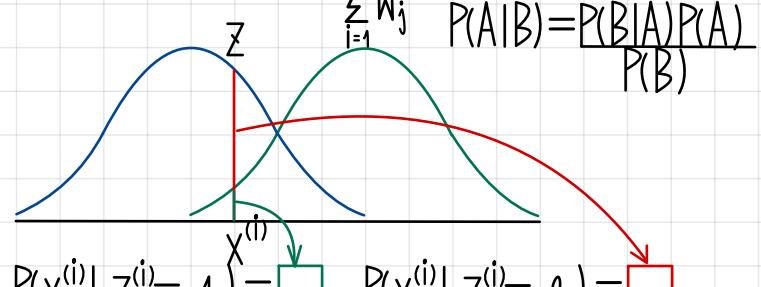
$w_j^{(i)} \leftarrow P(Z^{(i)}=j | X^{(i)}; \emptyset, \mu_1, \dots, \mu_k, \Sigma_1, \dots)$
↑ the responsibility of component j
for example i

M-Steps :

$$\begin{aligned}
 \emptyset_j &= \frac{1}{m} \sum_{i=1}^m w_j^{(i)} \\
 \mu_j &= \frac{\sum_{i=1}^m w_j^{(i)} X^{(i)}}{\sum_{i=1}^m w_j^{(i)}} \\
 \Sigma_j &= \frac{\sum_{i=1}^m w_j^{(i)} (X^{(i)} - \mu_j)(X^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}}
 \end{aligned}$$

optimal values of parameters under "soft" assignments $w_j^{(i)}$
 ↓ soft hard

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$



$$\boxed{} + \boxed{} P(Z^{(i)}=1 | X^{(i)}) = P(X^{(i)} | Z^{(i)}=1)P(Z^{(i)}=1)$$

$$\boxed{} + \boxed{} P(Z^{(i)}=2 | X^{(i)}) = P(X^{(i)} | Z^{(i)}=2)P(Z^{(i)}=2)$$

$$P(X^{(i)}) = P(X^{(i)} | Z^{(i)}=j)P(Z^{(i)}=j)$$

$$w_j^{(i)} = \sum_{l=1}^k P(X^{(i)}|Z^{(i)}=j; \mu_j, \Sigma_j) P(Z^{(i)}=j; \phi_j)$$

$$\sum_{l=1}^k P(X^{(i)}|Z^{(i)}=l; \mu_j, \Sigma_j) P(Z^{(i)}=l; \phi_l)$$

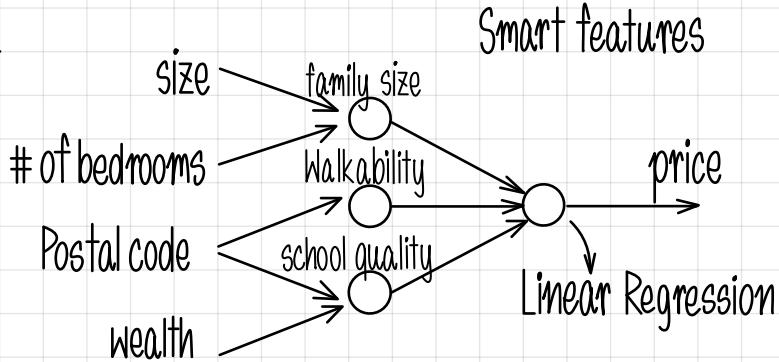
$$N(X^{(i)}; \mu_j, \Sigma_j)$$

EM is guaranteed to monotonically increase

$$l(\phi, \mu_1, \dots, \mu_k, \dots)$$

But can converge to local maximum

or become degenerate ($\Sigma_j \rightarrow 0$)



Deep Learning (a.k.a. neural networks)

$$f: X \rightarrow Y$$

$$\mathbb{R} \rightarrow \mathbb{R}^+$$

$$f(x) = \max(ax+b, 0)$$

ReLU : Rectified Linear Unit



It is possible to model **any** function to arbitrary accuracy with a sufficiently complex network of such unit

Many types of neural networks are **universal approximates**

(Note : "output" layer may not use the same activation function as all others)

Example : Predict housing prices

Based on: Postal code, size, # of bedrooms, wealth of neighborhood

In data science applications, we often combine raw feature into "Smart features"

Postal code \rightarrow Walkability
Wealth \rightarrow School quality

Loan underwriting : Demographic info, credit info
Predict : will they pay?

Logistic Regression : works well if payers & non-payers are linearly separable

Smart variable/feature : debt-to-income ratio

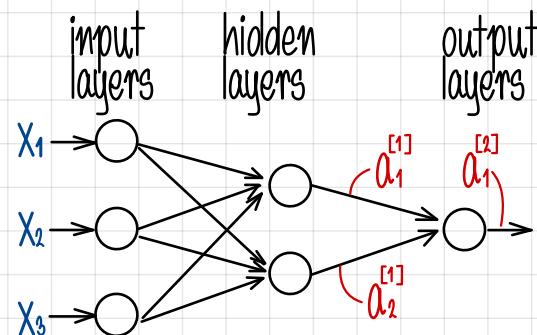
Deep neural networks perform end-to-end learning, from raw data input to output

Notation :

input feature : $x_1 + x_2 + \dots + x_n$

units in layers :

$a_i^{[j]}$ for the activation of the i^{th} units in the j^{th} layers



If Classification

$$a_1^{[2]} = \frac{1}{1+e^{-w^T x + b}}$$

Fully connected : all unit in one layer connect to all unit in next layer

$$x_1 = a_1^{[0]} \quad x_2 = a_2^{[0]} \quad \dots$$

Activation - 2 steps :

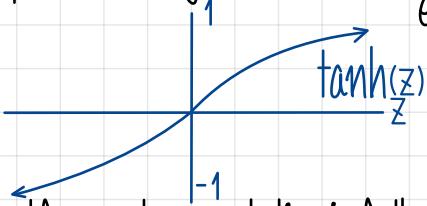
1. compute **Linear** activation $z = w^T x + b$

2. compute activation function $a = \sigma(z)$

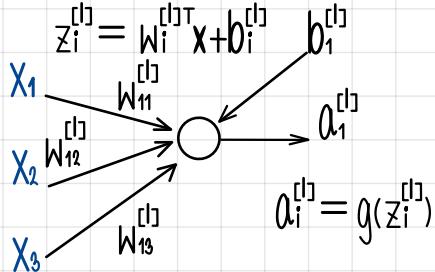
$$\text{output of a Classifier} : \delta(z) = \frac{1}{1+e^{-z}}$$

$$\text{Relu (mostly for hidden units)} : \max(z, 0)$$

$$\text{Hyperbolic tangent} : g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



Feedforward computation in fully connected model



Efficiency Calculation :

Try to make everything a vector or matrix operation (most efficiency calculation)

Numpy

No! Do not do this!

```
for l in range(L):
    z = np.zeros((N[l], 1))
    for i in range(N[l]):
        z[i] = b[i]
        for j in range(N[l-1]):
            z[i] = z[i] + w[i, j] * a[l-1, j]
            a[l, i] = g[z[i]]
```

Basic Linear
Algebra
Subroutines
Library

Do this! one called the "bias" library

```
for l in range(L):
    z[1] = b[1] + w[1] * a[1-1]
    a[1] = g[z[1]]
```

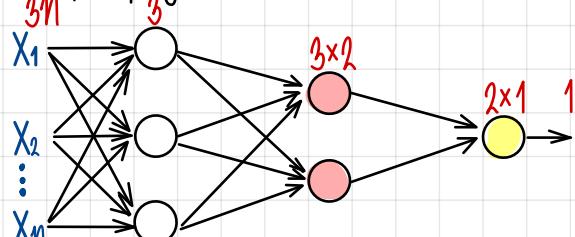
Inside $g(z)$:

matrix vector

```
def g(z):
    return 1/(1+np.exp(-z))
```

vector

Backpropagation



2 hidden Logistic sigmoid output

parameters : $3n+14$ (weight, bias)

initialization

$$w_{jk}^{[i]} \sim N(0, 0.1)$$

or Xavier initialization

$$w_{jk}^{[i]} \sim N(0, \sqrt{\frac{2}{n^{[i]} + n^{[i-1]}}})$$

Variance of activation is preserved across layers

Assume Relu for hidden units

Bias : small positive random or fixed values.

Output loss function for single example
(Stochastic Gradient Descent)

$$L(\hat{y}, y) = -[(1-y)\log(1-\hat{y}) + y\log\hat{y}]$$

For all parameters, we want

$$\begin{aligned} w^{[l]} &\leftarrow w^{[l]} - \alpha \frac{\partial L}{\partial w^{[l]}} \\ b^{[l]} &\leftarrow b^{[l]} - \alpha \frac{\partial L}{\partial b^{[l]}} \end{aligned}$$

Learning rate \downarrow Gradient

This in layer 3 :

$$\begin{aligned} \frac{\partial L}{\partial w^{[3]}} &= \frac{\partial [-(1-y)\log(1-\hat{y}) - y\log\hat{y}]}{\partial w^{[3]}} \\ &= -(1-y) \frac{\partial}{\partial w^{[3]}} \log(1-g(w^{[3]}a^{[2]}+b^{[3]})) \\ &\quad - y \frac{\partial}{\partial w^{[3]}} \log(g(w^{[3]}a^{[2]}+b^{[3]})) \\ &= \frac{-(1-y)(-1)g'(w^{[3]}a^{[2]}+b^{[3]})a^{[2]}}{1-g(w^{[3]}a^{[2]}+b^{[3]})} - \frac{y(g'(w^{[3]}a^{[2]}+b^{[3]})a^{[2]})}{g(w^{[3]}a^{[2]}+b^{[3]})} \end{aligned}$$

$g(z) : \text{sigmoid } \delta(z)$

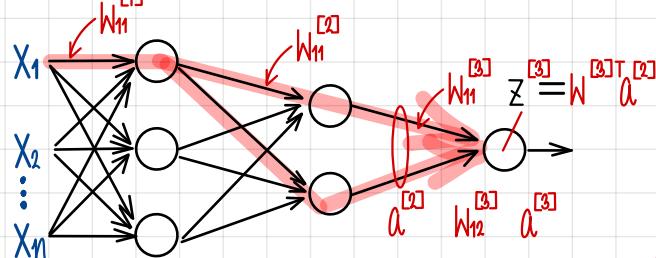
$$\delta(z)' = \delta(z)(1-\delta(z))$$

$$= (1-y)\delta'(w^{[3]}a^{[2]}+b^{[3]})a^{[2]} - y(1-\delta(w^{[3]}a^{[2]}+b^{[3]}))a^{[2]}$$

$$\frac{\partial L}{\partial w^{[3]}} = (a^{[2]} - y)a^{[2]}$$

Remember from Logistic Regression : $(\hat{y} - y)x$?

What about layer 2?



Chain Rule : $f(z)$ where $z = g(z)$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial x}$$

$a^{[3]}$ is a function of $z^{[3]}$
 $z^{[3]}$ is a function of $a^{[2]}$
 $a^{[2]}$ is a function of $z^{[2]}$
 $z^{[2]}$ depends on $W^{[2]}$

$$\frac{\partial L}{\partial W_{ij}^{[2]}} = \frac{\partial L}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}}$$

$$\begin{aligned} \frac{\partial L}{\partial W^{[2]}} &= \frac{\partial L}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial W^{[2]}} \\ &= (a^{[3]} - y) a^{[2]} \end{aligned}$$

Reused terms : "Delta" $\delta^{[3]}, \delta^{[2]}, \dots$

$$\delta^{[3]} = \frac{\partial L}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial W^{[2]}} \quad \frac{\partial L}{\partial W^{[2]}} = \delta^{[3]} a^{[2]}$$

$$\delta^{[2]} = \delta^{[3]} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \quad \frac{\partial L}{\partial W^{[1]}} = \delta^{[2]} a^{[1]}$$

Chain Rule for an effect on an output through multiple paths :

$$y = f(z) \quad u = g(z)$$

$$\begin{aligned} \frac{\partial y}{\partial x_i} &= \sum_j \frac{\partial y}{\partial u_j} \frac{\partial u_j}{\partial x_i} \\ \frac{\partial L}{\partial W_{ij}^{[2]}} &= \sum_k \frac{\partial L}{\partial a_k} \frac{\partial a_k}{\partial W_{ij}^{[2]}} \\ &= \sum_k \frac{\partial L}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a_k} \frac{\partial a_k}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a_j} \frac{\partial a_j}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial W_{ij}^{[2]}} \end{aligned}$$

Forward propagation

FC BP (SGD): Single Pattern x, y

given x, y

$a^{[0]} \leftarrow x$

for $l \in 1 \dots L$:

linear activation

$z^{[l]} \leftarrow W^{[l]} a^{[l-1]} + b^{[l]}$

nonlinear activation

$a^{[l]} \leftarrow \sigma^{[l]}(z^{[l]})$

$\# a^{[l]}$ is $\hat{y}!$

$\delta^{[l]} \leftarrow \frac{\partial L}{\partial z^{[l]}} \# a^{[l]} - y$ for Logistic sigmoid

Linear

Relu

Sigmoid

Tanh

More exotic ones!

for $l \in 1 \dots L$:

$\frac{\partial L}{\partial a^{[l]}} = \delta^{[l]} a^{[l]}$

$\frac{\partial L}{\partial W^{[l]}} = \delta^{[l]}$

$\frac{\partial L}{\partial b^{[l]}} = \delta^{[l]}$

if $l > 1$:

$$\delta^{[l-1]} \leftarrow \text{diag}(g^{[l-1]}(z^{[l-1]})) \cdot W^{[l]} \delta^{[l]}$$

What about Batch Gradient Descent

More accurate than SGD

More slow (in terms of data use) than SGD

Get best of Both worlds: Mini Batch Gradient Descent

1. sample a subset of training set

2. calculate Gradients over the whole Mini Batch

3. take a step

Iteration : one step for a Mini Batch

Epoch : one pass through entire training set

Example : training set size 1000, Batch size 10

we say "100 Iteration per Epoch"

Momentum : Hyperparameter in SGD

(Along with learning rate, Batch size)

Aiming to reduce the noise due to sampling

optimum

Calculate a "velocity"

Instantaneous
Stochastic
Gradient

$$V_{dw}^{[l]} \leftarrow \beta V_{dw}^{[l]} - (1-\beta) \frac{\partial J}{\partial w^{[l]}}$$

Moving Average of the Stochastic steps

$$w^{[l]} \leftarrow w^{[l-1]} - \alpha V_{dw}^{[l]}$$

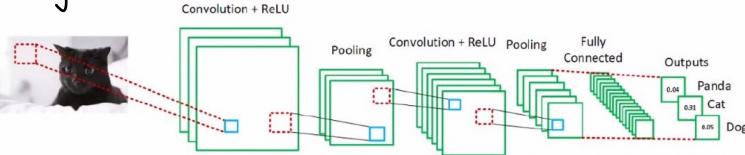
Smoothed Gradient

ADAM : Adaptive Sophisticated Adaptive Momentum

Optimize become another Hyperparameter

Convolutional Neural Networks (CNN)

Layers of a CNN



Reinforcement Learning

Supervised Learning

data driven → pair every input with target

Reinforcement Learning

data driven → reward, delayed
more "hand-off"

→ autonomous robot
→ game playing → Atari 2600 → video input
→ joystick output
→ reward score
David Silver → Go
→ online multiplayer games
Starcraft

Protein folding prediction ← paper

Markov decision process

A tuple $(S, A, \{P_{sa}\}, \gamma, R)$

S : Set of states

A : Set of actions (things the agent can do)

$\{P_{sa}\}$: States transition probabilities

For each states $s \in S$, each action $a \in A$,

P_{sa} gives a distribution over S

$P_{saj}(s_k)$: probabilities that when taking action

a_j in state s_i that the government moves to state s_k

$\gamma \in [0,1]$ discount factor : how important is it to get reward early rather than late

$R : S \times A \rightarrow \mathbb{R}$ reward function

MDP (Markov decision process) :

Start in some state S_0

Agent chooses action $a_0 \in A$

Environment samples $S_1 \sim P_{s,a_0}$

[continue choosing actions, samplings states]

Total payoff : $R(S_0, a_0) + \gamma R(S_1, a_1) + \gamma^2 R(S_2, a_2) + \dots$

What is the best action to take in each state?

Policy is a function $\Pi : S \mapsto A$

Execute the policy by taking, in States s , to action $a = \Pi(s)$ (deterministic or randomized)

A policy Π has a value function $V^\Pi(s)$

$$V^\Pi(s) = \mathbb{E} [R(S_0, a_0) + \gamma R(S_1, a_1) + \gamma^2 R(S_2, a_2) + \dots | S_0 = s; \Pi]$$

↑ the average of the expression over
an infinite number of trials

Famous recursive relationship : Bellman equations

$$V^\Pi(s) = R(s) + \gamma \sum_{s' \in S} P_{\Pi(s)}(s') V^\Pi(s')$$

immediate
reward

$$V^\Pi(s) = R(s) + \gamma \mathbb{E} [R(s_1) + \gamma R(s_2) + \dots | \Pi]$$

If $|S| < \infty$, write down $V^\Pi(s)$ for each $s \in S$, R, P_{sa}, Π fixed & known, $|S|$ Linear equation in $|S|$ unknowns

4	5	6
3		
2	1	0

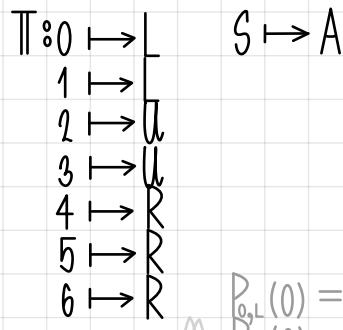
$S = \{0, 1, 2, 3, 4, 5, 6\}$ $R = 0 \mapsto 0$
 $A = \{L, R, U, D\}$ $1 \mapsto 0$
 \vdots $5 \mapsto 0$
 $\gamma = 0.9$ $6 \mapsto 100$

$$P_{SA} = 0, R \mapsto [1, 0, \dots, 0] \quad 1, U/D \mapsto [0, 1, 0, \dots, 0]$$

$$0, U \mapsto [1, 0, \dots, 0] \quad 1, L \mapsto [0, 0, 1, 0, \dots, 0]$$

$$0, D \mapsto [1, 0, \dots, 0] \quad 1, R \mapsto [0, 0, 0, 1, 0, \dots, 0]$$

$$0, L \mapsto [0, 1, 0, 0, \dots, 0]$$



$$V^\pi(0) = 0 + 0.9 \cdot 0.1 V^\pi(0) + 0.9 \cdot 0.9 V^\pi(1)$$

$$V^\pi(1) = 0 + 0.9 \cdot 0.1 V^\pi(1) + 0.9 \cdot 0.9 V^\pi(2)$$

$$V^\pi(2) = 0 + 0.9 \cdot 0.1 V^\pi(2) + 0.9 \cdot 0.9 V^\pi(3)$$

$$V^\pi(3) = 0 + 0.9 \cdot 0.1 V^\pi(3) + 0.9 \cdot 0.9 V^\pi(4)$$

$$V^\pi(4) = 0 + 0.9 \cdot 0.1 V^\pi(4) + 0.9 \cdot 0.9 V^\pi(5)$$

$$V^\pi(5) = 0 + 0.9 \cdot 0.1 V^\pi(5) + 0.9 \cdot 0.9 V^\pi(6)$$

$$V^\pi(6) = 100$$

$$0.91 V^\pi(0) - 0.91 V^\pi(1) = 0$$

$$0.91 V^\pi(1) - 0.91 V^\pi(2) = 0$$

$$0.91 V^\pi(2) - 0.91 V^\pi(3) = 0$$

$$0.91 V^\pi(3) - 0.91 V^\pi(4) = 0$$

$$0.91 V^\pi(4) + 0.91 V^\pi(5) = 0$$

$$0.91 V^\pi(5) + 0.91 V^\pi(6) = 0$$

$$V^\pi(6) = 1$$

$$Ax = b \quad x = A^{-1}b$$

$$A = \begin{bmatrix} 0.91 & -0.91 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0.91 & -0.91 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.91 & -0.91 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.91 & -0.91 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.91 & 0.91 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.91 & 0.91 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$x = \begin{bmatrix} V^\pi(0) \\ \vdots \\ V^\pi(6) \end{bmatrix}$$

Optimal value function $V^*(S)$

$$V^*(S) = \max_\pi V^\pi(S) \quad \text{expected future reward}$$

$$V^*(S) = R(S) + \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s') \quad \text{immediate reward}$$

Optimal policy $\Pi^*: S \mapsto A$

$$\Pi^*(S) = \arg\max_a \sum_{s' \in S} P_{sa}(s') V^*(s')$$

Value Iteration : $V^*(S) = V^{\Pi^*}(S) \geq V(S)$

For each state S estimate, update iteratively
 $V(S) \leftarrow 0$

Repeat until convergence

for each state S

$$V(S) \leftarrow R(S) + \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V(s')$$

Theorem : V converges to V^*

Policy Iteration :

Initialize Π randomly

Repeat until convergence

$$V(S) \leftarrow V^\pi(S) \quad (\text{solve linear system})$$

for each state S

$$\Pi(S) \leftarrow \arg\max_a \sum_{s' \in S} P_{sa}(s') V(s')$$

Theorem : Π converges to Π^* , V converges to V^*

So far : RL is easy

When # of state is small, # of action is small,

P_{sa} is given, $R(S)$ is given, state is observable

$R(S)$ is immediately observable

In real world problems, P_{sa} is not given

Two ways to deal with unknown environment :

- Model-Based methods estimate P_{sa}
via trial and error

- Model-Free methods avoid using P_{sa}
in the algorithm

If world is discrete state & discrete action, just use a table to keep track of how many times are a_p up to each state s when executing action a in state s

To explore, we typically use a simply randomized algorithm called ϵ -greedy to decide what to do

ϵ -greedy policy:

with probability ϵ , take a uniform random action

with probability $1-\epsilon$, take the current estimate best action

Introduce Q-Learning

Deal with infinite/partial state observation

Maximization bias, double Q-Learning

TD Methods (Temporal-Difference)

vs Monte Carlo

TD Prediction

similar to policy evaluation estimate V^π

Monte carlo RL method wait until real return are observed:

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

actual return obtained in state S_t following π to end of episode

TD Methods approximate this by performing updates

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

noisy

Tabular TD(0) for V^π

Input policy π , LR α

Initialize $V(S)$ for all $s \in S$ arbitrarily except $V(\text{terminal}) = 0$

For each step in episode:

$$a_t \leftarrow \pi(s_t)$$

$$s' \sim P_{sa}$$

$$V(s) \leftarrow V(s) + \alpha [R(s') + \gamma V(s') - V(s)]$$

Convergence is guaranteed if α is small and decrease overtime

We want a model free method to find π^*

when P_{sa} is unknown

Control Problem

SARSA, on-policy to policy iteration method instead of estimating V , we estimate Q (action value function):

$q_\pi(s, a)$: expected discounted return after executing action a in state s , following π

$$Q(S_t, a_t) \leftarrow Q(S_t, a_t) + \alpha [R(S_{t+1}) + \gamma Q(S_{t+1}, a_{t+1}) - Q(S_t, a_t)]$$
$$(S_t, a_t, r_t, S_{t+1}, a_{t+1})$$

TD Methods: SARSA Updated Equation

Tabular Methods $Q(S^t, a^t)$: Table lookup

What about infinite state/action spaces?

partial observability
infinite

discrete continuous ranges

output a parameter distribution p_{sa}

SARSA

$$Q(S_t, a_t) \leftarrow Q(S_t, a_t) + \alpha [R(S_{t+1}) + \gamma Q(S_{t+1}, a_{t+1}) - Q(S_t, a_t)]$$
$$(S_t, a_t, r_t, S_{t+1}, a_{t+1})$$

We must execute a_t, a_{t+1} before updating

Q-Learning

$$Q(S_t, a_t) \leftarrow Q(S_t, a_t) + \alpha [R(S_{t+1}) + \max_a Q(S_{t+1}, a) - Q(S_t, a_t)]$$

ϵ -greedy
 $p(1-\epsilon)$ $p(\epsilon)$

follow Q pick random action

exploration will select action according to current Q function

$a_{t+1} \neq a$

V^*

Q-Learning is "off-policy" - learns the Optimal policy rather than the Optimal ϵ -greedy policy

$q_\pi^*(S, \pi)$

Maximization Bias

Q-Learning: select a using $\max_a Q(S_{t+1}, a)$

the Maximization leads to bias

consider problem with a state s for which

$Q(s, a)$ is truly 0 for all a

However, due to noise/random initial values,

$\max_a Q(s, a)$ will always be positive

$\therefore V(s)$ is always overestimated

Double Q-Learning

$\forall s, a$, initialize $Q_1(s, a), Q_2(s, a)$ arbitrarily, except $Q(\text{terminal}, \cdot) = 0$

For each episode:

Initialize s

For each step in the episode:

choose action a , observing r, s' with probability $1/2$:

$$Q_1(s, a) \leftarrow Q_1(s, a) + \alpha [R + \gamma Q_2(s', \arg\max_a Q_1(s, a)) - Q_1(s, a)]$$

else:

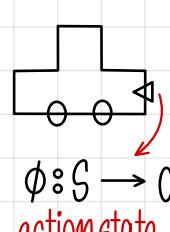
$$Q_2(s, a) \leftarrow Q_2(s, a) + \alpha [R + \gamma Q_1(s', \arg\max_a Q_2(s, a)) - Q_2(s, a)]$$

Eliminates Maximization Bias

$p_{sa} = \frac{\# \text{ of times we took action } a \text{ in state } s \text{ and got to } s'}{\# \text{ of times we took action } a \text{ in state } s}$

Dealing with infinite & partially observable state spaces

Example: self-driven vehicle



size of space: infinite
environment/state is only
partially observable via
ovr sensors
use sensor measurements as
a proxy for the state

Assume there exists an observation function Q
that transduces the world to a set of sensor
measurements with noise

Example sensor: monocular image sequence
sample preprocessing to make observation useful

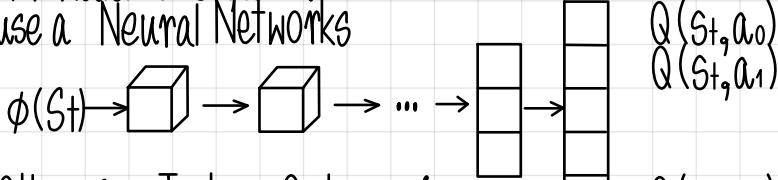
Mnih et al. (2013) (Deepmind / Google)
Atari Games (ALE / Gym)

frames 1 2 3 4 5 ...
colors 210x160 images

representation of state s_t
is a stack of 4 grayscale
frames

COMPACT REPRESENTATION OF ACTUAL STATE
(OPTICAL FLOW CAN GIVE INFO ABOUT TRAJECTORY /
VELOCITY). (28,224 elements!)
↓ 256 values per element

How model $Q(s_t, a_t)$?
use a Neural Networks



Otherwise, Just use Q-Learning
(actually, Double Q-Learning or a variant of π)

Replay Buffer: Save

$s_t, a_t, r_t, a_{t+1}, \text{done}$
Tuples in a Replay Buffer. Q-Learning is done
in batches sampled from the Buffer