

Lecture 3. Perceptron Learning Algorithm (PLA)

Efficiently find a perfect classifier for linearly separable dataset.

Notation:

$$\text{Decision Rule: } w_0x_0 + w_1x_1 + \dots + w_dx_d \stackrel{y=+1}{\geq} b \\ -b + \sum_{i=1}^d w_i x_i \stackrel{y=-1}{\geq} 0$$

Let $w_0 = -b$, $x_0 = 1$,

Augmented vectors:

$$\underline{w} = (w_0 = -b, w_1, \dots, w_d) \in \mathbb{R}^{d+1}$$

$$\underline{x} = (x_0 = 1, x_1, \dots, x_d) \in \mathbb{R}^{d+1}$$

$$\text{New decision rule: } \underline{w}^T \cdot \underline{x} \stackrel{y=+1}{\geq} 0$$

$$y = h_w(\underline{x}) \triangleq \text{sign}(\underline{w}^T \underline{x}) . \text{sign}(t) = \begin{cases} +1, & t > 0 \\ -1, & t < 0 \end{cases}$$

Suppose training set D is linearly separable, find $\underline{w} \in \mathbb{R}^{d+1}$ s.t.

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N I(y_n \neq h_w(x_n)) = 0$$

Input: Training set D that is linearly separable.

Output: $\underline{w} \in \mathbb{R}^{d+1}$ that achieves $E_{in}(\underline{w}) = 0$.

Initialization: arbitrary, e.g. $\underline{w} = \underline{0}$.

Step 1: Check if $E_{in}(\underline{w}) = 0$.

If yes, stop & output current \underline{w} .

Step 2: Let (x_n, y_n) be a miss-classified example.

(including points on boundary).

If $y_n = +1$, $\underline{w} \leftarrow \underline{w} + \underline{x}_n$

If $y_n = -1$, $\underline{w} \leftarrow \underline{w} - \underline{x}_n$

Go to step 1

Lec 4

Binary linear classification

Decision Rule: $y = \text{sign}(\underline{w}^T \underline{x})$.

augmented: $\underline{w} = (w_0, w_1, \dots, w_d)$

$$\underline{x} = (x_0 = 1, x_1, \dots, x_d)$$

Perceptron learning Algorithm

Update step: Pick mis-classified example (x_n, y_n) , Note: $y_n = +1$ or -1

$$\underline{w} \leftarrow \underline{w} + y_n \underline{x}_n$$

Intuition

Possible real & predicted result.

\hat{y}_n	$w^T \underline{x}_n$ (pred.)	is correct prediction?	$\hat{y}_n w^T \underline{x}_n$
+1	+	✓	+
+1	-	✗	-
-1	+	✗	-
-1	-	✓	+
±1	0	✗	0

If (\underline{x}_n, y_n) is correctly classified, $\hat{y}_n w^T \underline{x}_n > 0$.

If (\underline{x}_n, y_n) is mis-classified, $\hat{y}_n w^T \underline{x}_n \leq 0$.

Updating: Suppose (\underline{x}_n, y_n) is mis-classified.

$$w_{\text{new}} = w + \hat{y}_n \underline{x}_n,$$

$$\hat{y}_n w_{\text{new}}^T \underline{x}_n = \hat{y}_n (w + \hat{y}_n \underline{x}_n)^T \underline{x}_n$$

$$= \hat{y}_n \cdot w^T \underline{x}_n + \hat{y}_n^2 \|\underline{x}_n\|^2$$

$> \hat{y}_n \cdot w^T \underline{x}_n$ (since we have a bias) $\rightarrow +1$

Strict improvement for \underline{x}_n !

But it could cause new mis-classification for other examples.

Theorem (Rosenblatt, 1957)

Given a linearly separable dataset, the PLA terminates in a finite number of steps yielding

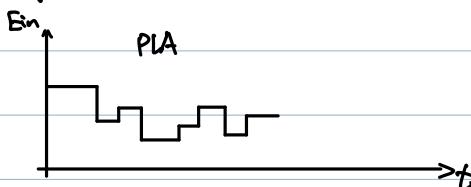
$$E_{in}(w) = 0.$$

Proof. Problem 13

Remark: The output from the algo. is not unique.

Pocket Algorithm

Extends PLA for datasets that are not linearly separable.



Keep the "best" weight vector upto iteration t in the pocket.

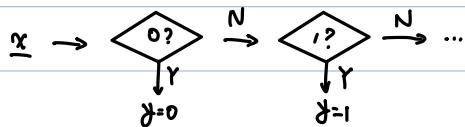
Algorithm:

- 0: Pick time horizon T .
- 1: Set pocketed weight vector \hat{w} to $w(0)$ of PLA.
- 2: For $t = 0, 1, 2, \dots, T-1$ do:
 - 3: Run PLA for one update to obtain $w(t+1)$.
 - 4: Evaluate $E_{in}(w(t+1))$
 - 5: If $E_{in}(w(t+1)) < E_{in}(\hat{w})$, set $\hat{w} = w(t+1)$.
 - 6: Return \hat{w} .

Heuristic algorithm, performs well in practice.

Note: Multiary classification can be implemented by a sequence of binary classification

e.g. 10 digits



Lecture 5.

Linear Regression

Recall: supervised learning

discrete y_n : classification

continuous y_n : regression

Training set:

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_m, y_m)\}$$

$$x_n \in \mathbb{R}^d, y_n \in \mathbb{R}$$

Decision Rule: (Hypothesis Set)

$$\hat{y} = h(x)$$

$$= w_0 + w_1 x_1 + \dots + w_d x_d$$

redefine augmented form: $\underline{w} = (w_0, w_1, \dots, w_d)$

$$\underline{x} = (x_0=1, x_1, \dots, x_d)$$

$$\hat{y} = h(\underline{x}) = \underline{w}^T \underline{x}$$

Criterion:

$$E_{in}(\underline{w}) = \underbrace{\frac{1}{N} \sum_{n=1}^N}_{e_n(\underline{w})} (\hat{y}_n - \underline{w}^T \underline{x}_n)^2$$

"average squared error"

$e_n(\underline{w})$ is the square error on the n th example.

Goal:

Given \mathcal{D} , find $\underline{w} \in \mathbb{R}^{d+1}$ to minimize $E_{in}(\underline{w})$

Example: x = advertising cost in one week ($d=1$)

y = sales in one week

Historical data

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$$

Fit a linear model

$$y = w_0 + w_1 x$$

$$\text{Refine model: } \underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \text{adv\$ on TV} \\ \text{on radio} \\ \text{on newspaper} \end{bmatrix}$$

$$(d=3)$$

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

larger $w_i \Rightarrow$ more profitable y_i .

More algebra:

1) Data matrix

$$X = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$$

2) Target vector (label)

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$$

3) Weight vector

$$\underline{w} \in \mathbb{R}^{d+1}$$

4) Model

$$\hat{\underline{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} \underline{w}^T \underline{x}_1 \\ \underline{w}^T \underline{x}_2 \\ \vdots \\ \underline{w}^T \underline{x}_N \end{bmatrix} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} \cdot \underline{w}$$

$$\hat{\underline{y}} = X \cdot \underline{w}$$

5) Error:

$$E_{in}(\underline{w}) = \frac{1}{N} \cdot \| \underline{y} - \hat{\underline{y}} \|^2$$

When is $E_{in}(\underline{w}) = 0$, i.e., $\underline{y} = \hat{\underline{y}}$?

$$y_n = \underline{w}^T \cdot \underline{x}_n, \quad n=1, 2, \dots, N.$$

of linear equation: N

of variable: $d+1$ (w_i)

In practice, $N \gg d$

\Rightarrow No exact solution for the eqn. Instead, we minimize $E_{in}(\underline{w})$.

Least Mean Squared Solution.

(to minimize $E_{in}(\underline{w}) = \frac{1}{N} \| \underline{y} - \hat{\underline{y}} \|^2$)

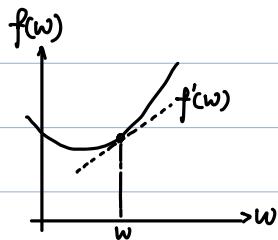
$$\text{Define } f(\underline{w}) = \| \underline{y} - \hat{\underline{y}} \|^2$$

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

$$= \sum_{n=1}^N (\underline{w}^T \underline{x}_n - y_n)^2$$

Def Gradient of $f(\underline{w})$

$$\nabla f(\underline{w}) = \begin{bmatrix} \frac{\partial f}{\partial w_0} \\ \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}$$



It points in the direction of steepest increase.

Claim: ① $\nabla f(\underline{w}) = 2 X^T (X \underline{w} - \underline{y})$

② Least squared solution \underline{w}_{LS} is st.

$$\nabla f(\underline{w}_{LS}) = 0$$

$$2 X^T (X \underline{w} - \underline{y}) = 0$$

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

Assumption:

The d+1 columns of $\underline{\mathbf{X}}$ are linearly independent.

\Leftrightarrow there exist at least d+1 rows of $\underline{\mathbf{X}}$ that are linearly independent.
 ↓
 data samples

Then $\text{rank}(\underline{\mathbf{X}}) = d+1$.

$\Leftrightarrow \underline{\mathbf{X}}^T \underline{\mathbf{X}}$ is invertible

$$\underline{\mathbf{w}}_{\text{LS}} = (\underline{\mathbf{X}}^T \underline{\mathbf{X}})^{-1} \underline{\mathbf{X}}^T \underline{\mathbf{y}} = \underline{\mathbf{X}}^+ \underline{\mathbf{y}}$$

$\underline{\mathbf{X}}^+ \triangleq (\underline{\mathbf{X}}^T \underline{\mathbf{X}})^{-1} \underline{\mathbf{X}}^T$ is called the pseudo-inverse of $\underline{\mathbf{X}}$.

$$\underline{\mathbf{X}}^+ \cdot \underline{\mathbf{X}} = \mathbb{I} \quad \text{but} \quad \underline{\mathbf{X}} \cdot \underline{\mathbf{X}}^+ \neq \mathbb{I}$$

Note: ① $\underline{\mathbf{y}} = \underline{\mathbf{X}} \underline{\mathbf{w}}$

$$\underline{\mathbf{w}} = ?$$

$$\textcircled{2} \quad \underline{\mathbf{y}}_{\text{LS}} = \underline{\mathbf{X}} \cdot \underline{\mathbf{w}}_{\text{LS}} = \underbrace{\underline{\mathbf{X}} \underline{\mathbf{X}}^+}_{\text{projection matrix}} \underline{\mathbf{y}}$$

Lecture 6

Geometric Interpretation of least squares

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

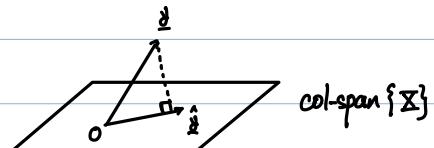
$$\text{LS solution: } \underline{\mathbf{w}}_{\text{LS}} = (\underline{\mathbf{X}}^T \underline{\mathbf{X}})^{-1} \underline{\mathbf{X}}^T \underline{\mathbf{y}}$$

$\hat{\underline{\mathbf{y}}} = \underline{\mathbf{X}} \underline{\mathbf{w}}$: linear combination of the columns of $\underline{\mathbf{X}}$,

$\Rightarrow \hat{\underline{\mathbf{y}}}$ is a vector in col-span of $\{\underline{\mathbf{x}}\}$

To minimize $\| \underline{\mathbf{y}} - \hat{\underline{\mathbf{y}}} \|^2$ (distance between $\underline{\mathbf{y}}$ & $\hat{\underline{\mathbf{y}}}$)

The best $\hat{\underline{\mathbf{y}}}$ (i.e., $\hat{\underline{\mathbf{y}}}_{\text{LS}}$) is the projection of $\underline{\mathbf{y}}$ onto col-span $\{\underline{\mathbf{x}}\}$.



\Rightarrow Every column of $\underline{\mathbf{X}}$ is orthogonal to $\underline{\mathbf{y}} - \hat{\underline{\mathbf{y}}}_{\text{LS}}$

$$\Rightarrow \underline{\mathbf{X}}^T (\underline{\mathbf{y}} - \hat{\underline{\mathbf{y}}}_{\text{LS}}) = 0$$

$$\Rightarrow \underline{\mathbf{X}}^T (\underline{\mathbf{y}} - \underline{\mathbf{X}} \underline{\mathbf{w}}_{\text{LS}}) = 0$$

$$\Rightarrow \underline{\mathbf{X}}^T \underline{\mathbf{X}} \underline{\mathbf{w}}_{\text{LS}} = \underline{\mathbf{X}}^T \underline{\mathbf{y}} \Rightarrow \underline{\mathbf{w}}_{\text{LS}} = (\underline{\mathbf{X}}^T \underline{\mathbf{X}})^{-1} \underline{\mathbf{X}}^T \underline{\mathbf{y}} \text{ as expected}$$

Regularized Linear Regression / Least Squares

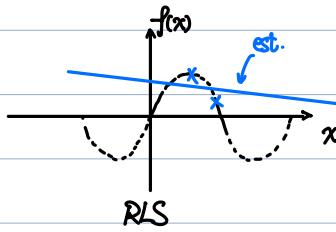
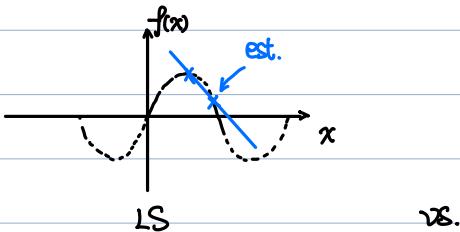
$$\min_{\underline{\mathbf{w}}} \| \underline{\mathbf{X}} \underline{\mathbf{w}} - \underline{\mathbf{y}} \|^2 + \lambda \|\underline{\mathbf{w}}\|^2$$

penalize large $\|\underline{\mathbf{w}}\|$

Motivation: avoid "overfitting" Ent-Ein

Textbook: Target function: $f(x) = \sin(\pi x)$

Training Set: 2 points (randomly sampled)



$$f(\underline{w}) = \|\underline{X}\underline{w} - \underline{y}\|^2 + \lambda \|\underline{w}\|^2$$

$$\frac{\partial f}{\partial w_k} = \frac{\partial}{\partial w_k} (\|\underline{X}\underline{w} - \underline{y}\|^2 + \lambda \sum_{j=0}^d w_j^2)$$

$$= 2[\underline{X}^T(\underline{X}\underline{w} - \underline{y})]_k + 2\lambda w_k$$

$$= 2[\underline{X}^T(\underline{X}\underline{w} - \underline{y}) + \lambda \underline{w}]_k \quad \text{only pick the } k^{\text{th}} \text{ element of the vector}$$

$$\nabla f(\underline{w}) = 2[\underline{X}^T(\underline{X}\underline{w} - \underline{y}) + \lambda \underline{w}]$$

$$\text{Want } \nabla f(\underline{w}) = \underline{0}$$

$$(\underline{X}^T \underline{X} + \lambda I) \cdot \underline{w}_{\text{RLS}} = \underline{X}^T \underline{y}$$

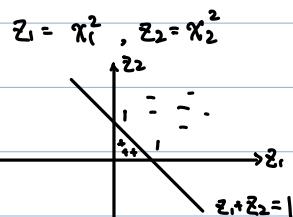
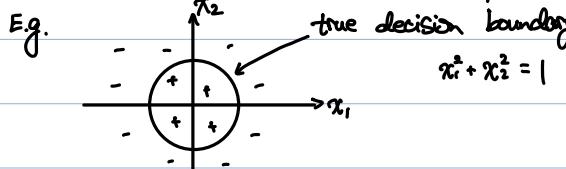
$$\underline{w}_{\text{RLS}} = (\underline{X}^T \underline{X} + \lambda I)^{-1} \cdot \underline{X}^T \underline{y}$$

$$\textcircled{1} \quad \underline{w}_{\text{RLS}} = \underline{w}_{\text{LS}} \text{ iff. } \lambda = 0$$

\textcircled{2} RLS side benefit more numerically stable.

Non-Linear Transformation

Linear function may not be sufficient.



$$\text{Suppose, } h(\underline{z}) = \text{sign}(z_1 + z_2 - 1)$$

$$\text{then } g(\underline{x}) = \text{sign}(x_1^2 + x_2^2 - 1)$$

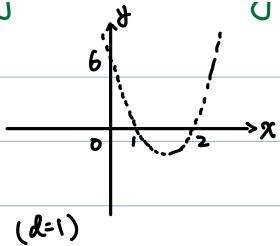
Let $\underline{z} = \Phi(\underline{x})$ be a non-linear transform.

Let $h(\underline{z}) = \underline{w}^T \underline{z}$ be a linear classifier in \underline{z} space.

Then $g(\underline{x}) = h(\Phi(\underline{x}))$ is a non-linear classifier in \underline{x} space.

$\Phi(\cdot)$ is called a feature transform.

E.g. Quadratic Regression



$$\underline{z} = (z_0=1, z_1=x, z_2=x^2).$$

$$y = \underline{w}^T \cdot \underline{z}$$

$$\begin{aligned} &= w_0 + w_1 z_1 + w_2 z_2 \\ &= w_0 + w_1 x + w_2 x^2 \end{aligned}$$

$$\underline{z}_1 = (1, 0, 0), \quad y_1 = 6$$

$$\underline{z}_2 = (1, 1, 1), \quad y_2 = 0$$

$$\underline{z}_3 = (1, 2, 4), \quad y_3 = 0$$

$$\Rightarrow \underline{z} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix}, \quad \underline{y} = \begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix}$$

$$\underline{w}_{LS} = (\underline{z}^T \underline{z})^{-1} \cdot \underline{z}^T \cdot \underline{y} = \underline{z}^{-1} \cdot \underline{y} = \begin{bmatrix} 6 \\ -9 \\ 3 \end{bmatrix}$$

$$y = 6 - 9x_1 + 3x_2$$

$$= 6 - 9x + 3x^2.$$

Lecture 7

Logistic Regression

Recap: Given $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$

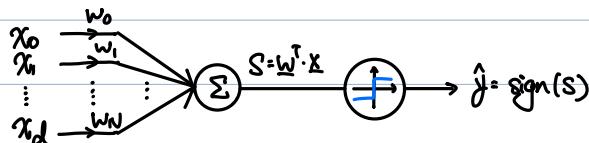
unknown target function $y = f(x)$,

Hypothesis: $\hat{y} = h(x)$ where $h \in H$.

Linear classification:

$$\underline{x} \in \mathbb{R}^{d+1}, \quad y \in \{+1, -1\}$$

$$\hat{y} = \text{sign}(\underline{w}^T \underline{x})$$

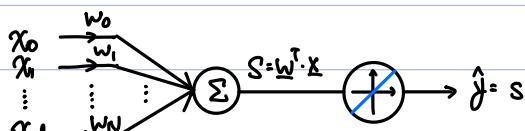


$$e(\underline{w}) = 1 \{ \hat{y} \neq \text{sign}(s) \}$$

Linear Regression

$$\underline{x} \in \mathbb{R}^{d+1}, \quad y \in \mathbb{R},$$

$$\hat{y} = \underline{w}^T \underline{x}$$



$$e(\underline{w}) = (\hat{y} - s)^2$$

Example:

What if we want randomness in f ?

E.g. \underline{x} = average (fat, sugar) in diet,

y = heart attack

Want to say "large $w^T \underline{x} \Rightarrow$ more likely that $y=+1$ "

Want new target function

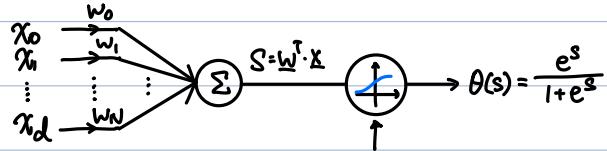
$$f_p(\underline{x}) = \Pr\{y=+1 | \underline{x}\}$$

"Soft decision"

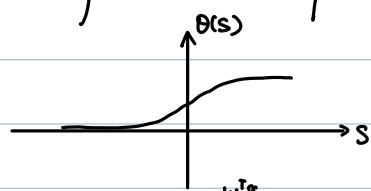
⇒ Solution:

Logistic Regression

$$\underline{x} \in \mathbb{R}^{d+1}, y \in \{+1, -1\}$$



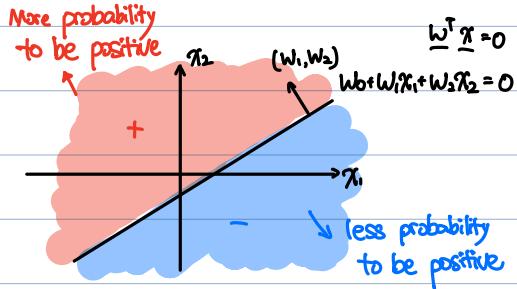
$\theta(s)$: Sigmoid function "S-shaped". $\tanh = \frac{e^s - e^{-s}}{e^s + e^{-s}}$



$$\text{Hypothesis: (guess)} \quad \Pr\{y=+1 | \underline{x}\} = \theta(s) = \theta(w^T \underline{x}) = \frac{e^{w^T \underline{x}}}{1 + e^{w^T \underline{x}}}$$

Note: ① Geometry:

($d=2$)



$$\textcircled{2} \quad \Pr\{y=-1 | \underline{x}\} = 1 - \Pr\{y=+1 | \underline{x}\}$$

$$= 1 - \frac{e^{w^T \underline{x}}}{1 + e^{w^T \underline{x}}}$$

$$= \frac{1}{1 + e^{w^T \underline{x}}}$$

$$= \frac{e^{-w^T \underline{x}}}{1 + e^{-w^T \underline{x}}}$$

$$= \theta(-w^T \underline{x}) = \theta(-s)$$

③ Notation:

$$\hat{P}_{\underline{w}}(1|\underline{x}) = \theta(\underline{w}^T \underline{x})$$

estimate for $\Pr\{y=+1|\underline{x}\}$

$$\hat{P}_{\underline{w}}(-1|\underline{x}) = \theta(-\underline{w}^T \underline{x})$$

$$\Rightarrow \hat{P}_{\underline{w}}(y|\underline{x}) = \theta(y \underline{w}^T \underline{x}) \quad . \quad y \in \{+1, -1\}$$

$$= \frac{e^{y \underline{w}^T \underline{x}}}{1 + e^{y \underline{w}^T \underline{x}}}$$

$$= \frac{1}{1 + e^{-y \underline{w}^T \underline{x}}}$$

Error Criterion:

For the n^{th} example in training data,

$$e_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n|\underline{x}_n) \leftarrow \text{"log-loss function"}$$

$$= \log(1 + e^{-y_n \underline{w}^T \underline{x}_n})$$

Note:

If $\underline{w}^T \underline{x} > 0$, and $y=+1 \Rightarrow \text{loss} = 0$.

$\underline{w}^T \underline{x} < 0$, and $y=-1 \Rightarrow \text{loss} = 0$.

Example: $\underline{x} \in \mathbb{R}^{d+1}$:

Suppose output of logistic regression is

$$\hat{P}_{\underline{w}}(-1|\underline{x}_n) = 0.8$$

$$\hat{P}_{\underline{w}}(+1|\underline{x}_n) = 0.2$$

If $y_n = +1$, then $e_n(\underline{w}) = -\log(0.2) = 1.39$

If $y_n = -1$, then $e_n(\underline{w}) = -\log(0.8) = 0.22$

Note: loss is smaller if $y_n = -1$, because our $\hat{P}_{\underline{w}}$ assigns higher probability for $y_n = -1$.

$$\text{Suppose } \hat{P}_{\underline{w}}(-1|\underline{x}_n) = 0.999$$

$$\hat{P}_{\underline{w}}(+1|\underline{x}_n) = 0.001,$$

If $y_n = +1$, $e_n(\underline{w}) = -\log(0.001) = 10$

If $y_n = -1$, $e_n(\underline{w}) = -\log(0.999) = 10^{-4}$

Logistic Regression

Given $\mathcal{D} = \{\underline{x}_1, y_1\}, \dots, \{\underline{x}_N, y_N\}$, find $\underline{w} \in \mathbb{R}^{d+1}$,

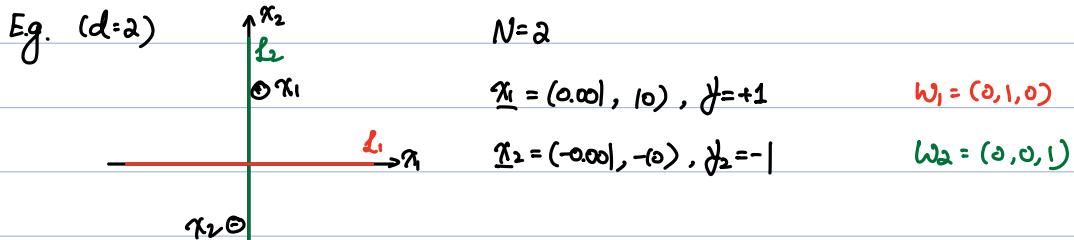
to minimize

$$E_n(\underline{w}) = \frac{1}{N} \sum_{n=1}^N e_n(\underline{w})$$

$$= \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n \underline{w}^T \underline{x}_n})$$

$$\begin{aligned} \text{Why use } \hat{E}_m(\underline{w}) &= -\log P_{\underline{w}}(\underline{d}_n | \underline{x}_n) ? \\ &= \log(1 + e^{-\underline{d}_n \underline{w}^T \underline{x}_n}) \end{aligned}$$

Benefit over linear classification



For linear.

$$\left. \begin{array}{l} E_m(w_1) = 0 \\ E_m(w_2) = 0 \end{array} \right\} \text{can't tell which is better}$$

For logistic regression

$$\begin{aligned} E_m(w_1) &= \frac{1}{2} (\log(1 + e^{-\underline{w}_1^T \underline{x}_1}) + \log(1 + e^{+\underline{w}_1^T \underline{x}_2})) \\ &= \frac{1}{2} (\log(1 + e^{-0.001}) + \log(1 + e^{-0.001})) \\ &= 0.693 \\ E_m(w_2) &= \frac{1}{2} (\log(1 + e^{-\underline{w}_2^T \underline{x}_1}) + \log(1 + e^{\underline{w}_2^T \underline{x}_2})) \\ &= \frac{1}{2} (\log(1 + e^{-10}) + \log(1 + e^{-10})) \\ &= 5 \times 10^{-5} \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} w_2 \text{ is better}$$

Note: What if $w_2 = (0, 0, 100)$?

Same line L_2 , but loss $E_m(w)$ is lower.

\Rightarrow Logistic regression typically requires regulation:

$$E_m(w) + \lambda \|w\|^2$$

Maximum likelihood Viewpoint

Let $\{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$ be the training dataset

$$\begin{aligned} \text{Consider } P(d_1, \dots, d_N | x_1, \dots, x_N) &\triangleq P\{ \text{1st label} = y_1, \dots | \text{1st example} = x_1, \dots \} \\ &= \prod_{n=1}^N P(d_n | x_n) \end{aligned}$$

by hypothesis

$$= \prod_{n=1}^N \hat{P}_{\underline{w}}(d_n | x_n).$$

Maximum-likelihood approach:

Find $\underline{w} \in \mathbb{R}^{d+1}$ that maximizes $P(d_1, \dots, d_N | x_1, \dots, x_N)$

$$\Leftrightarrow \text{maximize } \frac{1}{N} \log \prod_{n=1}^N \hat{P}_{\underline{w}}(d_n | x_n)$$

$$\begin{aligned}
 &= \frac{1}{N} \sum_{n=1}^N \log \hat{P}_w(\hat{y}_n | \underline{x}_n) \\
 \Leftrightarrow \text{minimize } &\frac{1}{N} \sum_{n=1}^N -\log \hat{P}_w(\hat{y}_n | \underline{x}_n) \\
 &= \frac{1}{N} \sum_{n=1}^N E_{in}(w) \\
 &= E_{in}(w)
 \end{aligned}$$

Cross-Entropy Viewpoint

$$E_{in}(w) = -\frac{1}{N} \sum_{n=1}^N [1(\hat{y}_n=+1) \log \hat{P}_w(+1 | \underline{x}_n) + 1(\hat{y}_n=-1) \log \hat{P}_w(-1 | \underline{x}_n)]$$

Define: Suppose P & Q are two probability distributions over $\mathcal{X} = \{\underline{x}_1, \dots, \underline{x}_N\}$.

Eg. Poisson with mean α , $\mathcal{X} = \{0, 1, 2, \dots\}$, $P(x) = \frac{\alpha^x}{x!} e^{-\alpha}$

The **cross-entropy** between P & Q (measurement of distance) is

$$CE(P, Q) = -\sum_{i=1}^m P(x_i) \log Q(x_i),$$

For n -th example, consider the distribution

$$\begin{aligned}
 P_n &= (\Pr\{\hat{y}_n=+1\}, \Pr\{\hat{y}_n=-1\}) \\
 &= \begin{cases} (1, 0), & \text{if } \hat{y}_n=1 \\ (0, 1), & \text{if } \hat{y}_n=-1 \end{cases} \\
 &= (1(\hat{y}_n=+1), 1(\hat{y}_n=-1))
 \end{aligned}$$

Let $Q_n = (\hat{P}_w(+1, \underline{x}_n), \hat{P}_w(-1, \underline{x}_n))$

(our estimate for the distribution of \hat{y}_n given \underline{x}_n).

$$\Rightarrow E_{in}(w) = \frac{1}{N} \sum_{n=1}^N CE(P_n, Q_n).$$

Lec 09

Gradient descent

Example: How to minimize $E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-\underline{x}_n w^\top \underline{x}_n})$

Given a differentiable function $f: \mathbb{R}^m \rightarrow \mathbb{R}$

want to $\min_{x \in \mathbb{R}^m} f(x)$

Gradient descent is a numerical approach.

$N=1$: (for illustration)

$x \in \mathbb{R}$

$$f'(x) \triangleq \frac{df}{dx}$$

if x at optimal, i.e. $x = x^*$, $f'(x) = 0$, [x^* : lowest point]

if $x > x^*$: $f'(x) > 0 \Rightarrow f(x)$ increases w.r.t. x

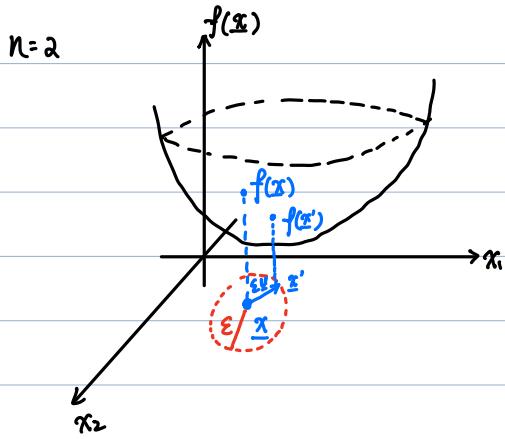
if $x < x^*$: $f'(x) < 0 \Rightarrow f(x)$ decreases w.r.t. x

1. Initialize $\underline{x} = \underline{x}_0$
 2. If $f'(\underline{x}) = 0$, then stop
 3. If $f'(\underline{x}) > 0$, then $\underline{x} = \underline{x} - \underline{\epsilon}$
 4. If $f'(\underline{x}) < 0$, then $\underline{x} = \underline{x} + \underline{\epsilon}$
 5. Go to step 2.
- $\underline{\epsilon}$: step size $\begin{cases} \text{large: less accurate} \\ \text{small: slow progress} \end{cases}$

Use $\underline{\epsilon}_k$, where k is the iteration number:

$$\text{e.g. } \underline{\epsilon}_k \sim \frac{1}{k}$$

General n :



Given current location \underline{x} , what should be the next step?

$$\underline{x}' \leftarrow \underline{x} + \underline{\epsilon} \underline{u}$$

$\underline{\epsilon}$ = step size

\underline{u} = direction of update. $\|\underline{u}\|=1$.

Idea: choose direction that maximizes $f(\underline{x}') - f(\underline{x})$

$$\begin{aligned} f(\underline{x}') &= f(\underline{x} + \underline{\epsilon} \underline{u}) \\ &\stackrel{\text{Taylor}}{=} f(\underline{x}) + \underline{\epsilon} \underline{u}^T \nabla f(\underline{x}) + O(\underline{\epsilon}^2) \\ \nabla f(\underline{x}) &= \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \end{aligned}$$

↑ negligible

Pick \underline{u} to minimize $\underline{u}^T \nabla f(\underline{x})$

$$\underline{u}^T \underline{v} = \|\underline{u}\| \cdot \|\underline{v}\| \cos \theta \geq -\|\underline{u}\| \|\underline{v}\| \text{ equity when } \cos \theta = -1.$$

$$\underline{u}^* = -\frac{\nabla f(\underline{x})}{\|\nabla f(\underline{x})\|}$$

Note: ① $\nabla f(\underline{x})$ points in the direction when $f(\underline{x})$ has max ascent.

② Magnitude $\|\nabla f(\underline{x})\|$ in the denominator can be absorbed into $\underline{\epsilon}$.

Gradient Descent Algorithm

Initialize \underline{x}_0 (typically pick at random)

For $t = 0, 1, 2, \dots$

1. Compute $\underline{g}_t = \nabla f(\underline{x}_t)$, $\nabla f(\underline{x}_t) = \nabla_{\underline{x}} f(\underline{x}) |_{\underline{x}=\underline{x}_t}$

2. Select direction $\underline{l}_t = -\underline{g}_t$

3. Update $\underline{x}_{t+1} = \underline{x}_t + \varepsilon_t \nabla f$

4. Go to step 1. until stopping criteria are reached

Note: ① ε_t : "learning rate"

② Condition for stopping: $\nabla f(\underline{x}_t) = 0$ is reasonable.

Lec 10

Gradient Descent

$$\underline{x}_{t+1} = \underline{x}_t - \varepsilon_t \nabla f(\underline{x}_t)$$

E.g. Logistic Regression:

$$\begin{aligned} \text{How to minimize } E_{\text{in}}(\underline{w}) &= \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-\underline{w}^T \underline{x}_n}) \\ &= \frac{1}{N} \sum_{n=1}^N e_n(\underline{w}) \end{aligned}$$

$$\text{GD: } \underline{w}_{k+1} = \underline{w}_k - \varepsilon_k \nabla E_{\text{in}}(\underline{w}_k)$$

$$= \underline{w}_k - \varepsilon_k \frac{1}{N} \sum_{n=1}^N \nabla e_n(\underline{w}_k)$$

$$\begin{aligned} \nabla e_n(\underline{w}_k) &= \nabla \underline{w} \log(1 + e^{-\underline{w}^T \underline{x}_n}) \\ &= \frac{1}{1 + e^{-\underline{w}^T \underline{x}_n}} \nabla \underline{w} (1 + e^{-\underline{w}^T \underline{x}_n}) \\ &= \frac{e^{-\underline{w}^T \underline{x}_n}}{1 + e^{-\underline{w}^T \underline{x}_n}} \nabla \underline{w} (-\underline{w}^T \underline{x}_n) \\ &= \frac{-\underline{w}^T \underline{x}_n}{1 + e^{-\underline{w}^T \underline{x}_n}} \end{aligned}$$

Stochastic Gradient Descent (SGD)

$$\text{GD: } \underline{w}_{k+1} = \underline{w}_k - \varepsilon_k \nabla E_{\text{in}}(\underline{w}_k)$$

SGD: At each iteration k , select uniformly random the n^{th} example from $\{1, 2, 3, \dots, N\}$

$$\text{and set } \underline{w}_{k+1} = \underline{w}_k - \varepsilon_k \nabla e_n(\underline{w}_k)$$

Note: ① $E[\nabla e_n(\underline{w}_k)]$

$$= \Pr(n=1) \cdot \nabla e_1(\underline{w}_k) + \dots + \Pr(n=N) \cdot \nabla e_N(\underline{w}_k)$$

$$= \frac{1}{N} \sum_{n=1}^N \nabla e_n(\underline{w}_k)$$

$$= \nabla E_{\text{in}}(\underline{w}_k)$$

i.e., $\nabla e_n(\underline{w}_k) = \text{mean} + \text{noise}$

$$= \nabla E_{\text{in}}(\underline{w}_k) + \text{noise}$$

a noisy estimate of the true gradient.

② Full GD ("batch" GD) has a complexity $O(Nd)$ per iteration

Stochastic GD is $O(d)$ Can do many more iterations than GD

(Often there is high redundancy in big datasets)

⇒ No need to consider all data points in each step.

M minibatch SGD

At each iteration, select M examples uniformly randomly:

$$\underline{x}_{(1)}, \underline{x}_{(2)}, \dots, \underline{x}_{(M)}$$

$$\underline{w}_{k+1} = \underline{w}_k - \frac{\epsilon_k}{M} \sum_{n=1}^M \nabla_{\underline{w}} e_n(\underline{w}_k)$$

- Better estimate of $\nabla E(\underline{w})$.

- Parallelization across multiple processors.

E.g. Logistic regression:

$$\nabla e_n(\underline{w}) = \frac{-\underline{y}_n \underline{x}_n}{1 + e^{\underline{y}_n \underline{w}^T \underline{x}_n}}$$

$$\text{SGD: } \underline{w}_{k+1} = \underline{w}_k + \epsilon_k \frac{\underline{y}_n \underline{x}_n}{1 + e^{\underline{y}_n \underline{w}^T \underline{x}_n}}$$

Note: close connection to the perceptron learning algo.

→ Suppose $(\underline{x}_n, \underline{y}_n)$ is mis-classified,

$$\underline{y}_n \underline{w}^T \underline{x}_n < 0,$$

$\Rightarrow e^{\underline{y}_n \underline{w}^T \underline{x}_n}$ is small.

SGD update rule is very close to perceptron, $\underline{w}_{k+1} \approx \underline{w}_k + \epsilon_k \underline{y}_n \underline{x}_n$.

→ Suppose $(\underline{x}_n, \underline{y}_n)$ is correctly classified:

$$\underline{y}_n \underline{w}^T \underline{x}_n > 0,$$

$\Rightarrow e^{\underline{y}_n \underline{w}^T \underline{x}_n}$ is large

SGD update rule $\Rightarrow \underline{w}_{k+1} \approx \underline{w}_k$

Lecture 11.

Multiclass Logistic Regression

Label $y \in \{1, 2, \dots, c\}$

Hypothesis:

Let $\Omega = \{\underline{w}(1), \dots, \underline{w}(c)\}$ be the weight vectors for the c classes.

Hypothesis that

$$\Pr\{\underline{y}_n = i | \underline{x}_n\} = \frac{e^{\underline{w}(i)^T \underline{x}_n}}{\sum_{j=1}^c e^{\underline{w}(j)^T \underline{x}_n}}, \quad i=1, 2, \dots, c. \quad \text{Softmax function}$$

$$\hat{p}(i | \underline{x}_n)$$

Error Criterion

$$e_n(\Omega) = -\log \hat{p}(y_n | \underline{x}_n) \quad \leftarrow \text{this is a vector}$$

$$\text{Note: } \nabla_{\underline{w}} e_n(\Omega) = \begin{bmatrix} \nabla_{\underline{w}(1)} e_n(\Omega) \\ \nabla_{\underline{w}(2)} e_n(\Omega) \\ \vdots \\ \nabla_{\underline{w}(c)} e_n(\Omega) \end{bmatrix}$$

SGD update:

For iteration k , $\underline{w}_k = \{\underline{w}_k(1), \dots, \underline{w}_k(c)\}$

Sample $n \sim \text{uniform } \{1, 2, \dots, N\}$.

For $i=1, 2, \dots, c$,

Compute $\nabla_{\underline{w}(i)} e_n(\underline{w})$

$$\underline{w}_{k+1}(i) = \underline{w}_k(i) - \varepsilon_k \nabla_{\underline{w}(i)} e_n(\underline{w}_k)$$

Compute $\nabla_{\underline{w}(i)} e_n(\underline{w})$: $\leftarrow n^{\text{th}}$ example (\underline{x}_n, y_n) is chosen

For $i=y_n$,

$$\nabla_{\underline{w}(i)} e_n(\underline{w}) = \nabla_{\underline{w}(i)} \left(-\log \frac{e^{\underline{w}(y_n)^T \underline{x}_n}}{\sum_{j=1}^c e^{\underline{w}(j)^T \underline{x}_n}} \right) \quad \leftarrow \underline{w}(y_n) \text{ is the weight of class } y_n$$

$$\begin{aligned} &= \nabla_{\underline{w}(i)} \left(-\underline{w}(y_n)^T \underline{x}_n + \log \left(\sum_{j=1}^c e^{\underline{w}(j)^T \underline{x}_n} \right) \right) \\ &= -\underline{x}_n + \frac{1}{\sum_{j=1}^c e^{\underline{w}(j)^T \underline{x}_n}} \cdot \nabla_{\underline{w}(i)} \left(\sum_{j=1}^c e^{\underline{w}(j)^T \underline{x}_n} \right) \end{aligned}$$

$$= -\underline{x}_n + \frac{e^{\underline{w}(i)^T \underline{x}_n}}{\sum_{j=1}^c e^{\underline{w}(j)^T \underline{x}_n}} \underline{x}_n.$$

For $i \neq y_n$,

$$\begin{aligned} \nabla_{\underline{w}(i)} e_n(\underline{w}) &= \nabla_{\underline{w}(i)} \left(-\underline{w}(y_n)^T \underline{x}_n + \log \left(\sum_{j=1}^c e^{\underline{w}(j)^T \underline{x}_n} \right) \right) \\ &= -\frac{e^{\underline{w}(i)^T \underline{x}_n}}{\sum_{j=1}^c e^{\underline{w}(j)^T \underline{x}_n}} \underline{x}_n \end{aligned}$$

Softmax Regression for $C=2$ (Binary classification)

$$\hat{P}(1|\underline{x}) = \frac{e^{\underline{w}\underline{x}}}{e^{\underline{w}\underline{x}} + e^{\underline{w}(2)^T \underline{x}}} = \frac{e^{(\underline{w}(1)-\underline{w}(2))^T \underline{x}}}{e^{(\underline{w}(1)-\underline{w}(2))^T \underline{x}} + 1}$$

$$\hat{P}(2|\underline{x}) = \frac{e^{\underline{w}(2)^T \underline{x}}}{e^{\underline{w}(1)^T \underline{x}} + e^{\underline{w}(2)^T \underline{x}}} = 1 - \hat{P}(1|\underline{x})$$

\equiv logistic regression with $\underline{w} = \underline{w}(1) - \underline{w}(2)$.

GD/SGD in linear regression

$$E_{\text{lin}}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N E_n(\underline{w})$$

$$E_n(\underline{w}) = (\underline{w}^T \underline{x}_n - y_n)^2$$

$$\begin{aligned} \nabla E_n(\underline{w}) &= \nabla_{\underline{w}} (\underline{w}^T \underline{x}_n - y_n)^2 \\ &= 2(\underline{w}^T \underline{x}_n - y_n) \cdot \nabla_{\underline{w}} (\underline{w}^T \underline{x}_n - y_n) \\ &= 2(\underline{w}^T \underline{x}_n - y_n) \cdot \underline{x}_n \end{aligned}$$

GD/SGD: as $k \rightarrow \infty$, \underline{w}_k converges to the least square solution

$$W_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

We have perfect solution, & why we need GD/SGD?

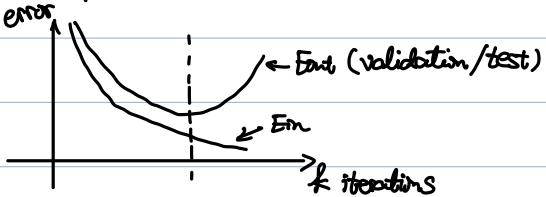
① Computational complexity:

GD: $O(Nd)$, SGD: $O(d)$, mini-batch SGD: $O(nd)$ per iteration

② Exact solution may not be desirable (over-fitting)

\uparrow more important We only care about E_{out} (test error), not E_{in} .

In practice, we run a few iterations, do validation & stop when error over the validation dataset is small.



Lecture 12

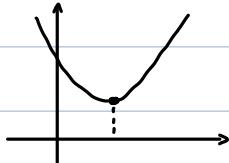
GD / SGD for non-convex Functions

Convex functions:

e.g. logistic regression

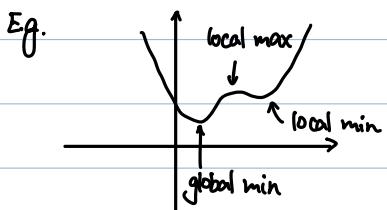
$$E_{\text{in}}(\mathbf{w})$$

linear regression



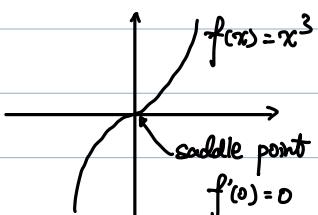
\Rightarrow global minimum. when $\nabla E_{\text{in}}(\mathbf{w}) = 0$

Non-convex functions:



Saddle points: $\nabla f(\mathbf{x}) = 0$, but \mathbf{x} is neither a local minimum nor a local maximum.

1-D:



Very slow progress near saddle points.

\Rightarrow pre-set # iteration runs out.

n-D space:

① 1-D saddle points for any direction.

② local minimum in some direction & local maximum in some other direction.

③ both 1 & 2 combined.

⇒ Highly likely to have saddle points in high-dimensional space.

Stopping Conditions:

to reduce the chances of stopping at a saddle points

① $\nabla E(\underline{w}) \approx 0$

② $E(\underline{w})$ is small.

③ # of iterations is large.

SGD with Momentum

Basic SGD:

$$\underline{g}_k = \nabla E(\underline{w}_k)$$

$$\underline{w}_{k+1} = \underline{w}_k - \epsilon_k \underline{g}_k$$

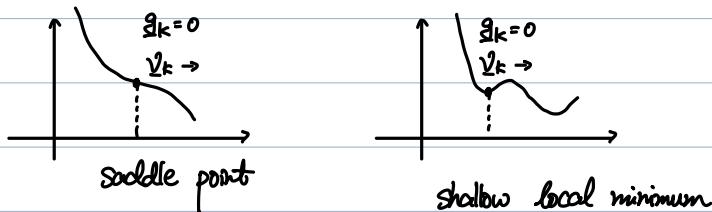
SGD + momentum

$$\underline{g}_k = \nabla E(\underline{w}_k)$$

$$\underline{v}_k = -\epsilon_k \underline{g}_k + \underline{m} \underline{v}_{k-1}$$

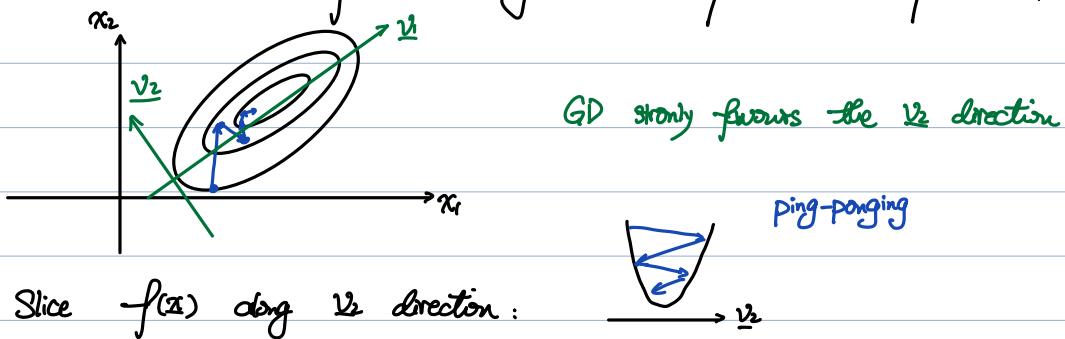
$$\underline{w}_{k+1} = \underline{w}_k + \underline{v}_k$$
 accumulation of previous movements

① Momentum helps SGD escape flat regions



"heavy ball momentum"

② Momentum can lead to faster convergence even for convex function.



ping-ponging

But overall trend ("momentum")

is moving along v_1 direction:

SGD + momentum! (less ping-pong)

Nesterov Momentum (83)

$$\dot{w}_k = -\varepsilon_k \nabla_{w_k} (\underbrace{w_k + \mu w_{k-1}}_{\text{intermediate point}}) + \mu \dot{w}_{k-1}$$

$$w_{k+1} = w_k + \dot{w}_k$$

Can prove better convergence.

For convex functions:

- Full GD: distance between w_k & w^* $\mathcal{O}(k)$
- with nesterov momentum: $\mathcal{O}(\frac{1}{k^2})$

Lecture 13

Multi-layer perception (MLP)

Recall perception: $\hat{y} = \text{sign}(w^T \cdot x)$

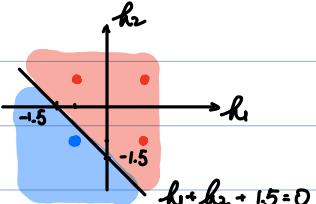
For non-separable dataset:

1. Pocket algorithm
2. Non-linear transformation \Rightarrow linear separable.

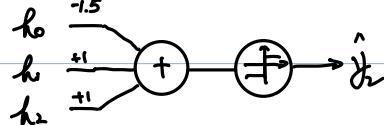
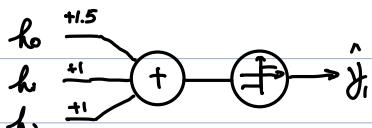
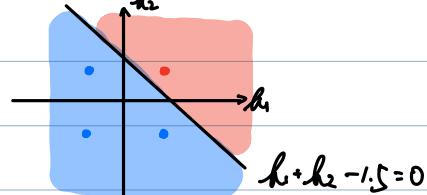
Note: too many candidates, not necessary a good transformation

Example:

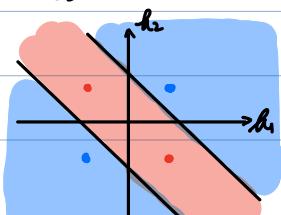
$$1. \hat{y}_1 = \text{OR}(h_1, h_2)$$



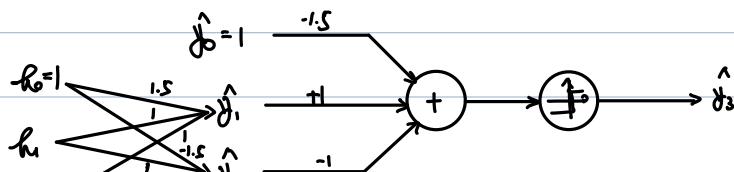
$$2. \hat{y}_2 = \text{AND}(-h_1, h_2)$$



$$3. \hat{y}_3 = \text{XOR}(h_1, h_2)$$



Multi-layer Perception (2-layer neural network)



$$\hat{y}_3 = \text{AND}(\hat{y}_1, \hat{y}_2)$$

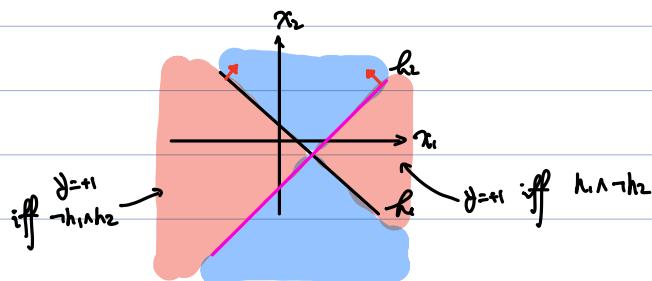
However, we would like to compute $y = f(x)$.

One more layer: function $h_1(x)$ & $h_2(x)$

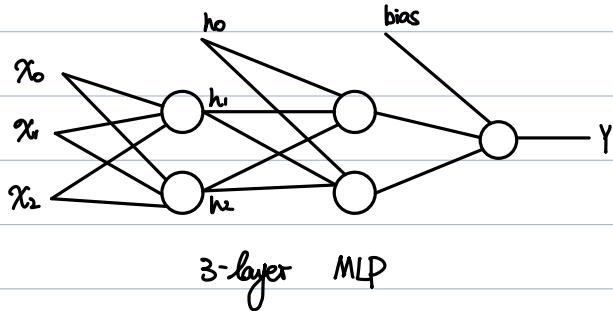
Example:

$$h_1(x) = \text{Sign}(x_1 + x_2 - 1)$$

$$h_2(x) = \text{Sign}(-x_1 + x_2 + 2)$$



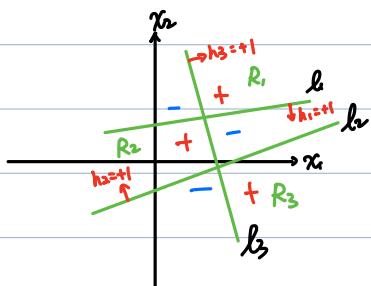
$$Y = \text{XOR}(h_1(x), h_2(x))$$



Lecture 14

Example:

Fact: Any decision rule can be expressed by intersection of hyperplanes
can be interpreted by using a 3-layer MLP.



$$R_1 : h_1 = -1, h_2 = +1, h_3 = +1 \iff \text{AND}(h_1, h_2, h_3)$$

$$R_2 : h_1 = +1, h_2 = +1, h_3 = -1 \iff \text{AND}(h_1, h_2, \bar{h}_3)$$

$$R_3 : h_1 = +1, h_2 = -1, h_3 = +1 \iff \text{AND}(h_1, \bar{h}_2, h_3)$$

$$y = +1, \text{ iff } x \in R_1 \cup R_2 \cup R_3$$

$$y = \text{OR}(R_1, R_2, R_3)$$

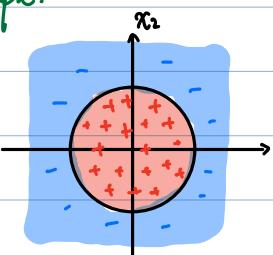
3-layer MLP:

layer 1: compute h_1, h_2, h_3

layer 2: AND function

layer 3: OR function

Example:



$$\text{Target: } y = \text{sign}(1 - x_1^2 - x_2^2)$$

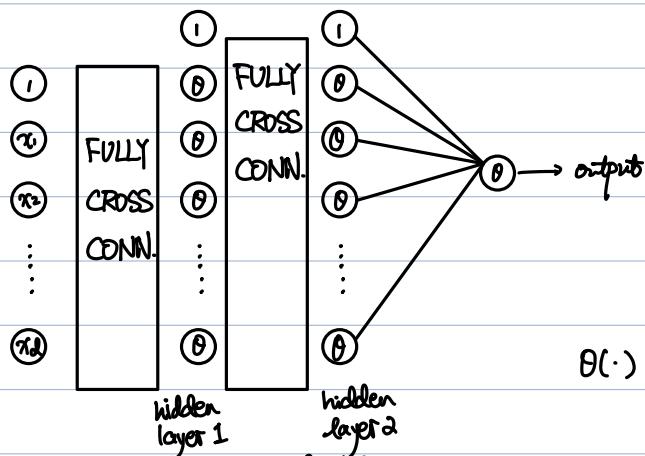
Can use more lines for approx.

Fact: A 3-layer MLP with sufficiently many hidden units in each layer can approximate any smooth decision boundary.

Neural Network (NN)

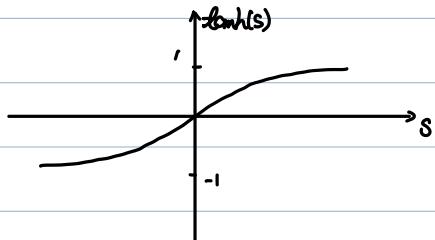


Usually use a softer function than $\text{sgn}(x)$ for easier optimization or learning.



$\theta(\cdot)$ is the activation function.

$$\text{Example: } \theta(s) = \tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$



*Universal Approx. Thm.: If the activation function is smooth, single layer can be sufficiently enough to approx. any function.

- $\theta(s) \approx s$ if $|s| \approx 0$
- Saturation if $|s| > 1$
- $\theta(s) = 1 - \theta(-s)$

Notation:

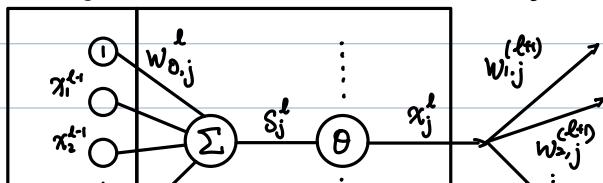
Input layer : $l=0$

Hidden layer : $1 \leq l \leq L-1$

Output layer : $l=L$

$d^{(l)}$: number of nodes in layer l : $0 \leq l \leq L$ (excluding bias)

$w_{i,j}^{(l)}$: weight connection node i in layer $l-1$ to node j in layer l





Input to node j in layer l : $s_j^{(l)}$

Output of node j in layer l : $x_j^{(l)}$

$$x_j^{(l)} = \theta(s_j^{(l)})$$

$$s_j^{(l)} = w_{0,j}^{(l)} + \sum_{i=1}^{d^{l-1}} w_{i,j}^{(l)} \cdot x_i^{(l-1)}$$

Lecture 15. ap.7.1 & 7.2

Vector notation

$$\underline{s}^{(l)} = \begin{bmatrix} s_1^{(l)} \\ s_2^{(l)} \\ \vdots \\ s_{d^l}^{(l)} \end{bmatrix}, \quad \theta(\underline{s}^{(l)}) = \begin{bmatrix} \theta(s_1^{(l)}) \\ \theta_2(s_2^{(l)}) \\ \vdots \\ \theta(s_{d^l}^{(l)}) \end{bmatrix}, \quad \underline{x}^{(l)} = \begin{bmatrix} 1 \\ x_1^{(l)} \\ x_2^{(l)} \\ \vdots \\ x_{d^l}^{(l)} \end{bmatrix} = \begin{bmatrix} 1 \\ \theta(\underline{s}^{(l)}) \end{bmatrix}$$

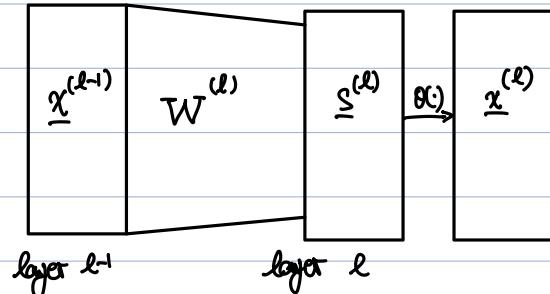
$W^{(l)} = \{w_{i,j}^{(l)}\}_{\substack{0 \leq i \leq d^{(l-1)} \\ 1 \leq j \leq d^{(l)}}}$ (matrix)

$$\underline{s}^{(l)} = (W^{(l)})^T \cdot \underline{x}^{(l-1)}$$

Forward Propagation

$$\underline{x}^{(0)} \xrightarrow{W^{(1)}} \underline{s}^{(1)} \xrightarrow{\theta(\cdot)} \underline{x}^{(1)} \xrightarrow{W^{(2)}} \underline{s}^{(2)}$$

$$\dots \longrightarrow \underline{x}^{(l-1)} \xrightarrow{W^{(l)}} \underline{s}^{(l)} \xrightarrow{\theta(\cdot)} \underline{x}^{(l)}$$



Pseudo Code

Input : $\underline{x}^{(0)} = (1, x_1, x_2, \dots, x_d)$

For $l = 1, 2, \dots, L$ do
 $\underline{s}^{(l)} = (W^{(l)})^T \cdot \underline{x}^{(l-1)}$
 $\underline{x}^{(l)} = [\underset{1}{\theta}(\underline{s}^{(l)})]$

Computational Complexity

of nodes: $V = \sum_{i=1}^L d^{(i)}$

of edges: $E = \sum_{l=0}^{L-1} (d^{(l)} + 1) d^{(l+1)}$

of computation: $O(V + E)$
 activation function Matrix Multiplication

Training of Neural Network

Loss function:

Let $\Omega = \{W^{(1)}, W^{(2)}, \dots, W^{(L)}\} \in \mathbb{R}^3$ be our model parameters

$l_n(\Omega)$ loss function w.r.t. example (x_n, y_n)

Example: linear regression for output layer:

square error: $(x_n^{(i)} - y_n)^2$

logistic: regression for output layer.

softmax: regression for output layer.

log-loss: $-\log \hat{P}(y_n | x_n)$

Stochastic Gradient Descent

In iteration k , suppose Ω_k is the current set of parameters.

- Select (x_n, y_n) at random, and compute.

$$\nabla_{\Omega} l_n(\Omega_k) = \frac{\partial l_n(\Omega_k)}{\partial W_{i,j}^{(l)}} , \forall i, j, l.$$

- $W_{i,j}^{(l)} \leftarrow W_{i,j}^{(l)} - \epsilon \cdot \frac{\partial l_n(\Omega_k)}{\partial W_{i,j}^{(l)}} , \forall i, j, l.$

Numerical Approach (finite difference)

Replace $W_{i,j}^{(l)}$ by $W_{i,j}^{(l)} + \Delta w$ while keeping all other weights fixed.

Let Ω_k^+ be the resultant set of parameters

$$\frac{\partial l_n(\Omega_k^+)}{\partial W_{i,j}^{(l)}} = \frac{l_n(\Omega_k^+) - l_n(\Omega_k)}{\Delta w}$$

$O(V+E)$ complexity to compute $l_n(\Omega_k^+)$ by forward propagation.

Overall complexity is $O(VE + E^2)$ for each iteration.

Typically, #of weights (Z) $\sim 10M$.

Lecture 16

Back Propagation

Note:

1. $e(\Omega)$ is memoryless. No knowledge about any previous layers.

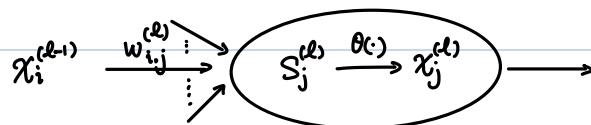
$$e(\Omega) = e(S^{(L)}, W^{(L)}, \dots, W^{(1)})$$

↑
all info up to layer l .

$$= e((S_i^{(L)}, \dots, S_j^{(L)}, \dots, S_d^{(L)}), W^{(L)}, \dots, W^{(1)})$$

2. Only $S_j^{(L)}$ depends on $w_{i,j}^{(L)}$.

$$\Rightarrow \frac{\partial e(\Omega)}{\partial w_{i,j}^{(L)}} = \frac{\partial e(\Omega)}{\partial S_j^{(L)}} \cdot \frac{\partial S_j^{(L)}}{\partial w_{i,j}^{(L)}}$$



$$S_j^{(l)} = w_{0,j}^{(l)} + \sum_{k=1}^{d^{(l-1)}} w_{k,j}^{(l)} \cdot x_k^{(l-1)}$$

$$\Rightarrow \frac{\partial S_i^{(l)}}{\partial w_{i,j}^{(l)}} = x_i^{(l-1)}$$

let $\delta_j^{(l)} \triangleq \frac{\partial e(\omega)}{\partial S_j^{(l)}}$ (Sensitivity of $e(\omega)$ to layer l inputs)

$$\Rightarrow \frac{\partial e(\omega)}{\partial w_{i,j}^{(l)}} = x_i^{(l-1)} \cdot \delta_j^{(l)}$$

$$\frac{\partial e(\omega)}{\partial w^{(l)}} \triangleq \left[\frac{\partial e(\omega)}{\partial w_{i,j}^{(l)}} \right] \Bigg|_{\begin{array}{c} 0 \leq i \leq d^{(l-1)} \\ 1 \leq j \leq d^{(l)} \end{array}}, \quad x^{(l-1)} = \begin{bmatrix} x_0^{(l-1)} \\ \vdots \\ x_{d^{(l-1)}}^{(l-1)} \end{bmatrix}, \quad \delta^{(l)} = \begin{bmatrix} \delta_1^{(l)} \\ \vdots \\ \delta_{d^{(l)}}^{(l)} \end{bmatrix}$$

$$\text{Then } \frac{\partial e(\omega)}{\partial w^{(l)}} = \underbrace{x^{(l-1)}}_{n \times 1} \underbrace{(\delta^{(l)})^T}_{1 \times d}$$

To compute $\delta^{(l)}$ "backwards"

① Consider $l = l$,

let $d^{(l)} = 1$ for simplicity.

$$S^{(l)} \rightarrow \Theta \rightarrow x^{(l)}, \quad g(x^{(l)}, y) = e(\omega)$$

$$\begin{aligned} \delta^{(l)} &= \frac{\partial e(\omega)}{\partial S^{(l)}} = \frac{\partial g(x^{(l)}, y)}{\partial S^{(l)}} \\ &= \frac{\partial g(x^{(l)}, y)}{\partial x^{(l)}} \cdot \frac{\partial x^{(l)}}{\partial S^{(l)}} \\ &= \frac{\partial g(x^{(l)}, y)}{\partial x^{(l)}} \cdot \theta'(S^{(l)}) \end{aligned}$$

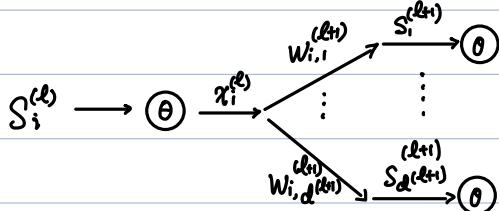
Example : Regression :

$$g(x^{(l)}, y) = (x^{(l)} - y)^2$$

$$\delta^{(l)} = 2(x^{(l)} - y) \cdot \theta'(S^{(l)})$$

② Intermediate node: $\delta_i^{(l)}$

$$\begin{aligned} \delta_i^{(l)} &= \frac{\partial e(\omega)}{\partial S_i^{(l)}} \\ &= \frac{\partial e(\omega)}{\partial x_i^{(l)}} \cdot \frac{\partial x_i^{(l)}}{\partial S_i^{(l)}} \\ &= \frac{\partial e(\omega)}{\partial x_i^{(l)}} \cdot \theta'(S_i^{(l)}) \end{aligned}$$



Recall: chain rule

$$\text{If } z = f(x, y) \text{, and } \begin{cases} x = g(s, t) \\ y = h(s, t) \end{cases}$$

$$\text{then } \frac{\partial z}{\partial s} = \frac{\partial z}{\partial x} \cdot \frac{\partial x}{\partial s} + \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial s}$$

$$\frac{\partial z}{\partial t} = \frac{\partial z}{\partial x} \cdot \frac{\partial x}{\partial t} + \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial t}$$

$$\Rightarrow \frac{\partial e(\omega)}{\partial x_i^{(l)}} = \sum_{j=1}^{d^{(l+1)}} \frac{\partial e(\omega)}{\partial S_j^{(l+1)}} \cdot \frac{\partial S_j^{(l+1)}}{\partial x_i^{(l)}} = \sum_{j=1}^{d^{(l+1)}} \delta_j^{(l+1)} \cdot w_{i,j}^{(l+1)}$$

$$\Rightarrow \delta_i^{(l)} = \left[\sum_{j=1}^{d^{(l+1)}} \delta_j^{(l+1)} \cdot w_{i,j}^{(l+1)} \right] \cdot \theta'(S_i^{(l)})$$

$$\Rightarrow \underline{\delta}^{(l)} = [\hat{W}^{(l+1)} \cdot \underline{\delta}^{(l+1)}] \otimes \theta'(S^{(l)})$$

↑
pointwise multiplication

where $\underline{\delta}^{(l)} = \begin{bmatrix} \delta_1^{(l)} \\ \vdots \\ \delta_d^{(l)} \end{bmatrix}, \theta'(S^{(l)}) = \begin{bmatrix} \theta'(S_1^{(l)}) \\ \vdots \\ \theta'(S_d^{(l)}) \end{bmatrix}$

$$\hat{W}^{(l+1)} \triangleq \begin{bmatrix} w_{1,1}^{(l+1)} & \cdots & w_{1,d}^{(l+1)} \\ \vdots & \ddots & \vdots \\ w_{d,1}^{(l+1)} & \cdots & w_{d,d}^{(l+1)} \end{bmatrix}$$

Lecture 17

Backpropagation Algorithm

Input: $(x, y), \Omega = \{W^{(1)}, W^{(2)}, \dots, W^{(L)}\}$

Output: $\frac{\partial e(\Omega)}{\partial w_{ij}^{(l)}}, l=1, 2, \dots, L$

1. Run forward propagation to compute:

$$\{S^{(l)}, x^{(l)}\}, l=1, 2, \dots, L$$

$$e(\Omega) = g(x^{(1)}, y)$$

$$2. \underline{\delta}^{(1)} = [\nabla_{x^{(1)}} g(x^{(1)}, y)] \otimes \theta'(S^{(1)})$$

3. For $l=L-1 \rightarrow 1$ do:

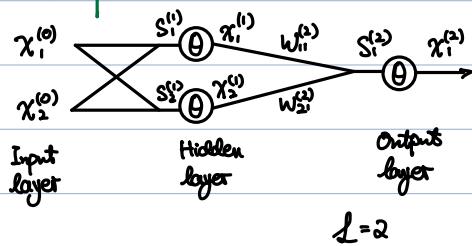
$$\underline{\delta}^{(l)} = [\hat{W}^{(l+1)} \cdot \underline{\delta}^{(l+1)}] \otimes \theta'(S^{(l)})$$

$$\frac{\partial e(\Omega)}{\partial w_{ij}^{(l)}} = x^{(l-1)} (\underline{\delta}^{(l)})^T$$

End.

Complexity: $O(E) \leftarrow$ One forward propagation to compute all $\frac{\partial e}{\partial w_{ij}}$

For example:



$$e(\Omega) = (x_i^{(2)} - y)^2$$

layer 2: $\underline{\delta}_i^{(2)} = \frac{\partial e}{\partial x_i^{(2)}} \cdot \frac{\partial x_i^{(2)}}{\partial S_i^{(2)}}$

$$= 2(x_i^{(2)} - y) \cdot \theta'(S_i^{(2)})$$

$$\frac{\partial e}{\partial w_{ii}^{(2)}} = x_i^{(1)} \cdot \underline{\delta}_i^{(2)}, \frac{\partial e}{\partial w_{21}^{(2)}} = x_2^{(1)} \cdot \underline{\delta}_1^{(2)}$$

$$\text{Layer 1: } \delta_1^{(1)} = \theta'(S_1^{(0)}) \cdot \delta_1^{(2)} \cdot w_{11}^{(2)}$$

$$\delta_2^{(1)} = \theta'(S_2^{(0)}) \cdot \delta_2^{(2)} \cdot w_{21}^{(2)}$$

$$\frac{\partial e}{\partial w_{ij}^{(1)}} = x_i^{(0)} \cdot \delta_j^{(1)}, \quad i, j \in \{1, 2\}$$

Note:

① Can do more:

$$\begin{aligned} \frac{\partial e}{\partial x_i^{(0)}} &\triangleq \delta_i^{(0)} = \frac{\partial e}{\partial S_i^{(0)}} \\ &= \frac{\partial e}{\partial S_1^{(0)}} \cdot \frac{\partial S_1^{(0)}}{\partial x_1^{(0)}} + \frac{\partial e}{\partial S_2^{(0)}} \cdot \frac{\partial S_2^{(0)}}{\partial x_1^{(0)}} \\ &= \delta_1^{(1)} \cdot w_{11}^{(1)} + \delta_2^{(1)} \cdot w_{21}^{(1)} \end{aligned}$$

Similarly, we have:

$$\frac{\partial e}{\partial x_2^{(0)}} = \delta_1^{(1)} \cdot w_{12}^{(1)} + \delta_2^{(1)} \cdot w_{22}^{(1)}$$

Note: 1. if too sensitive to the input, system may change drastic with small changes in inputs.
Too sensitive to noise, overfit.

2. what if we force

$$w_{12}^{(1)} = w_{21}^{(1)} \triangleq w ?$$

$$\begin{array}{c} x_1^{(0)} \xrightarrow[w]{\diagdown} S_1^{(1)} \Theta \\ x_2^{(0)} \xrightarrow[w]{\diagup} S_2^{(1)} \Theta \end{array}$$

$$\begin{aligned} \frac{\partial e}{\partial w} &= \frac{\partial e}{\partial S_1^{(1)}} \cdot \frac{\partial S_1^{(1)}}{\partial w} + \frac{\partial e}{\partial S_2^{(1)}} \cdot \frac{\partial S_2^{(1)}}{\partial w} \\ &= \delta_1^{(1)} \cdot x_2^{(0)} + \delta_2^{(1)} \cdot x_1^{(0)} \end{aligned}$$

More general case CNN.

Implementation of Neural Network

1) Choice of activation function

→ Sigmoid

$$\theta(s) = \frac{e^s}{1+e^s}$$

→ Tanh

$$\theta(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$

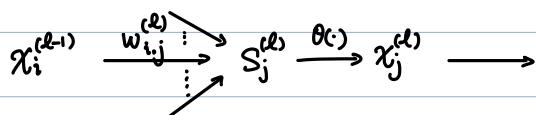
→ Rectified Linear Unit $\theta(s) = \max(0, s)$

- Simple calculation ($\theta'(s)=0, \forall s < 0 \Rightarrow$ Dead neuron \Rightarrow Less calculation)

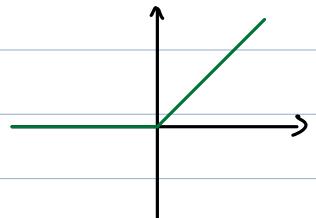
- Avoids non-linearity effect

(default in modern NNs)

Issue: Dead neuron



$$S_j^{(l)} < 0 \Rightarrow \{ x_j^{(l)} = 0 \}$$



$$\left| \theta'(S_j^{(l)}) = 0 \right.$$

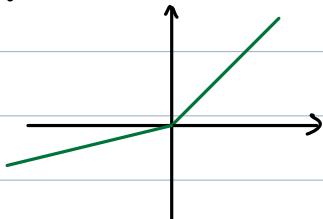
$$\Rightarrow \frac{\partial e}{\partial w_{ij}^{(l)}} = 0, \forall i$$

- Weights no longer updated.

- Cause ~40% neurons to die during training.

→ leaky ReLU

→ parametric ReLU



α is selected during training.

$$\theta(s) = \begin{cases} s, & s > 0 \\ \alpha s, & s \leq 0 \end{cases}$$

e.g. $\alpha=0.1$

Lecture 18

Implementation of Neural Network (Cont.)

2) Input pre-processing

① Normalization

$$D = \{(x_1, y_1), \dots, (x_N, y_N)\}$$

e.g. $\underline{x} \begin{cases} \text{salary} \\ \# \text{ of children} \end{cases}$

$$\text{Sample mean: } \underline{\mu} = \frac{1}{N} \sum_n x_n$$

$$\text{Sample standard deviation} \quad \underline{\sigma} = \sqrt{\frac{1}{N-1} \sum_n (x_n - \underline{\mu})^2} \quad \leftarrow \text{pointwise}$$

$$\underline{x}_n \leftarrow \frac{x_n - \underline{\mu}}{\underline{\sigma}} \quad \text{pointwise}$$

⇒ Each element has zero mean & unit variance

② Data augmentation

$$(\underline{x}, y) \leftarrow (F(\underline{x}), y)$$

$F(\cdot)$ can be rotation, flip, crop, scaling, etc.

Eg. Principle Component Analysis

3) Weight Initialization (sec 8.4)

- Small weight

works ok with small NN.

Vanishing gradient in deep NN.

(weights do not transfer info between layers)

- Large weight

signal grows over layers.

e.g. saturation of $\sigma(x)$.

E.g. Xavier Initialization

0 for bias terms

$w_{i,j}^{(l)} \sim N(0, \frac{1}{n})$, mean of 0, variance of $\frac{1}{n}$

where $n = \# \text{ of edges from } l-1 \text{ to Node } j \text{ of layer } l$.

$$S_j^{(l)} = w_{0,j}^{(l)} + \sum_{i=1}^n x_i^{(l-1)} \cdot w_{i,j}^{(l)}$$

$$\text{var}(S_j^{(l)}) = 0 + n(\frac{1}{n}) = 1, \text{ (assume independence among } x_i\text{'s)}$$

Other choices:

$$\sim N(0, \frac{\frac{2}{3}}{\#\text{ of in} + \#\text{ of out}})$$

~ uniform distr. with var $\frac{1}{n}$

$$U(-\frac{\sqrt{6}}{n}, \frac{\sqrt{6}}{n})$$

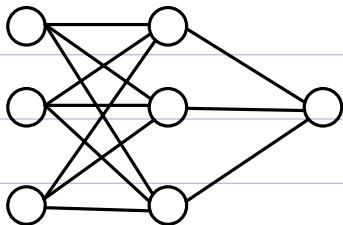
Note: With ReLU, need $\text{var}(w_{i,j}^{(l)}) \approx \frac{2}{n}$

4) Dropout

During training, in each SGD update, for each node with probability p set its output to zero in the forward path.

\Rightarrow During back propagation, only the weights connected to the active nodes are updated.

$$\frac{\partial e}{\partial w_{i,j}^{(l)}} = x_i^{(l-1)} \cdot \delta_j^{(l)}$$



Handwaving: - by multiple NNs at the same time.

Testing: Use the entire network, but scale the weights by the factor of $(1-p)$

$$w_{i,j}^{(l)} \leftarrow w_{i,j}^{(l)} \cdot (1-p)$$

5) Variation of SGD.

① Basic SGD

$$\underline{w}_{t+1} = \underline{w}_t - \varepsilon_t \nabla f(\underline{w}_t)$$

$\nabla f(\cdot)$ is a vector

② SGD with momentum

$$\underline{v}_t = \alpha \underline{v}_{t-1} - \varepsilon_t \nabla f(\underline{w}_t)$$

$$\underline{w}_{t+1} = \underline{w}_t + \underline{v}_t$$

③ Ada Gradient

Vary learning rate across different dimension.

$$\underline{c}_t = \underline{c}_{t-1} + \nabla f(\underline{w}_t) \otimes \nabla f(\underline{w}_t)$$

$$\underline{w}_{t+1} = \underline{w}_t - \varepsilon_t \nabla f(\underline{w}_t) \otimes \frac{1}{\sqrt{\underline{c}_t + 10^{-5}}}$$

↑
small number for stability.

Advantage:

smaller step size along steeper dimension.

Disadvantage:

Cumulate meaningless past

④ RMS prop

- Gradually forget distant history

$$\underline{c}_t = \alpha \underline{c}_{t-1} + (1-\alpha) \cdot \nabla f(\underline{w}_t) \otimes \nabla f(\underline{w}_t), \quad 0 < \alpha < 1$$

⑤ RMS-Prop + momentum

⑥ Adam

⑤ + correction factor

Lecture 19 Unsupervised Learning

$$D = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n\}, \quad \underline{x}_i \in R^d$$

- No labels

Eg. D: set of documents

Goal: group by topic

\underline{x}_i : histogram of word length in document i.

① Clustering

② Density estimation

③ Dimensionality reduction

Clustering (chapter 6.3.3)

partition D into K disjoint clusters s.t. the elements in each cluster are close to each other.

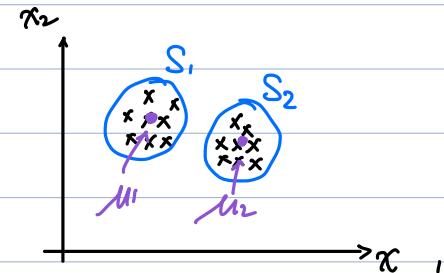
Given D , output

① Clusters S_1, S_2, \dots, S_k . $S_i \subseteq D$

$$\text{s.t. } S_i \cap S_j = \emptyset \quad \forall i \neq j$$

$$\bigcup_{i=1}^k S_i = D$$

② Cluster center: $\underline{m}_1, \underline{m}_2, \underline{m}_3, \dots, \underline{m}_k$, $\underline{m}_i \in \mathbb{R}^d$



Error Measure

$$E_j = \sum_{x_n \in S_j} \|x_n - \underline{m}_j\|^2. \text{ approximate error for cluster } S_j.$$

Given D and K , define

$$E_h(S_1, S_2, \dots, S_k, \underline{m}_1, \underline{m}_2, \dots, \underline{m}_k) = \sum_{j=1}^K E_j \\ = \sum_{n=1}^N \|x_n - \underline{m}(x_n)\|^2$$

where $\underline{m}(x_n) \equiv$ center of cluster for which x_n belongs

- Learning problem

$$\min_{S_1, S_2, \dots, S_K} E_h$$

- Applications

- Classification (e.g. documents, coins)

- Recommendation system

Optimal cluster is NP-hard!

Heuristic

K-means cluster

An alternating optimization approach with two subproblems

Subproblem #1.

Given S_1, \dots, S_k , find $\underline{m}_1, \dots, \underline{m}_k$ to minimize E_h .

$E_h = \sum_{j=1}^K E_j$, $E_j = \sum_{x_n \in S_j} \|x_n - \underline{m}_j\|^2$ depends only on S_j .

\Rightarrow Only need to consider

$$\min_{\underline{m}_j} E_j(\underline{m}_j) \iff \min_{\underline{m}_j} \sum_{x_n \in S_j} \|x_n - \underline{m}_j\|^2$$

$$\nabla E_j(\underline{m}_j) = 0$$

$$\Rightarrow \nabla_{\underline{m}_j} \sum_{x_n \in S_j} \|x_n - \underline{m}_j\|^2$$

$$= \sum_{x_n \in S_j} -2(x_n - \underline{m}_j) = 0$$

$$\Rightarrow \sum_{x_n \in S_j} x_n = \underline{m}_j |S_j|$$

$\underline{m}_j = \frac{1}{|S_j|} \sum_{x_n \in S_j} x_n$, sample average of points in S_j .

Subproblem #2

Given $\underline{m}_1, \dots, \underline{m}_k$, find S_1, \dots, S_k to minimize E_{in} .

\equiv Given $x_i \in D$, which cluster should it be associated?

$\underline{m}(x_i)$: cluster center for x_i .

$$E_{in} = \sum_{n=1}^N \|x_n - \underline{m}(x_n)\|^2$$

To minimize of $\|x_n - \underline{m}(x_n)\|$ depends only on x_n & $\underline{m}(x_n)$

$$\underline{m}(x_n) = \underset{\underline{m} \in \{\underline{m}_1, \dots, \underline{m}_k\}}{\operatorname{argmin}} \|x_n - \underline{m}\|^2$$

Assign x_n to the nearest cluster center.

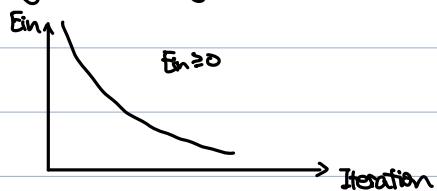
Lecture 20

K-means Algorithm

1. Initialize $\underline{m}_1, \dots, \underline{m}_k$ randomly
2. Construct S_1, \dots, S_k by solving subproblem #2 ("nearest cluster center")
3. Update $\underline{m}_1, \dots, \underline{m}_k$ by solving subproblem #1 ("centroid")
4. Repeat step 2 & 3 until converge.

Notes:

① Convergence is guaranteed

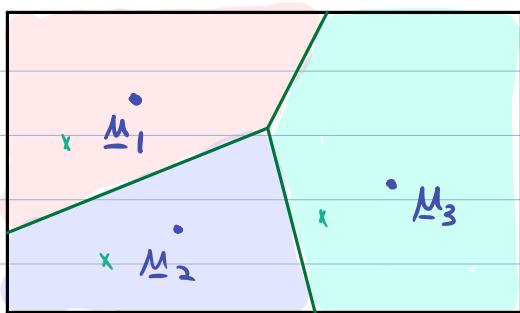


Monotone Convergence Theorem

② Locally optimal

Not global optimal in general

③ End Result:

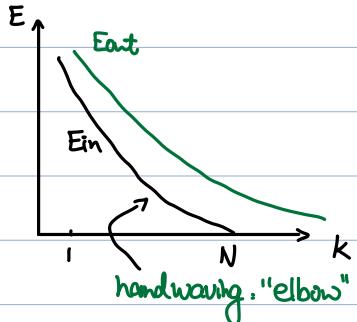


if $k=3$

k -means creates a voronoi diagram ("tessellation")

Test / inference: given new data point \underline{x} .

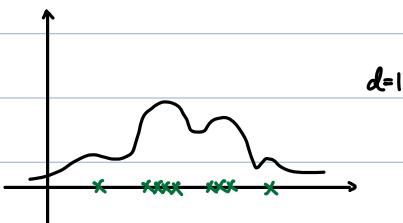
④ How to choose K ?



Density Estimation

$$D = \{\underline{x}_1, \dots, \underline{x}_N\}, \underline{x}_i \stackrel{iid}{\sim} p(\underline{x})$$

$$\text{output } \hat{p}(\underline{x}) \approx p(\underline{x})$$



Application: anomaly detection

Review:

Uniform Distribution

Discrete RV $X \sim \text{uniform in } \{1, 2, \dots, 10\}$

$$\text{PMF: } \Pr\{X=k\} = \frac{1}{10}, \quad 1 \leq k \leq 10$$

Continuous RV $X \sim \text{uniform in } (a, b)$

$$\Pr\{X = \frac{a+b}{2}\} = 0$$

$$\text{PDF: } f_X(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$$

Gaussian Distribution

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Question: Given sample points, how do we find the best line to fit density function?

1) Histogram Method

Assume $P(\cdot)$ is non-zero only over a bounded region cover region with uniform bins.

B_1, B_2, \dots, B_k , each with volume V . (Area / Range within the bin)

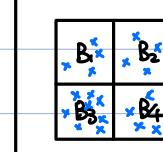
Let N_i be the # of samples in bin B_i

$$\text{Hypothesis: } \Pr\{X \in B_i\} = \frac{N_i}{N}$$

Assume that the distribution within each bin is uniform. \Rightarrow PDF: $\frac{1}{V}$

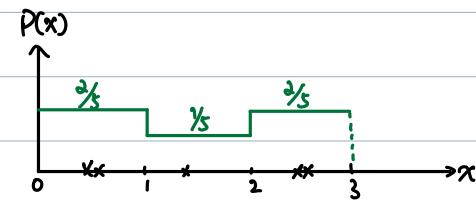
Total probability theorem

$$\hat{P}(x) = \sum_{i=1}^k \hat{P}(x | x \in B_i) \frac{N_i}{N}$$

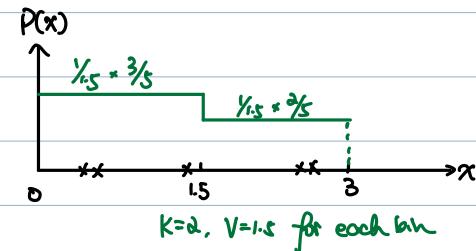


Since we assume uniform distribution in each bin,

$$\begin{aligned} \hat{P}(x) &= \sum_{i=1}^k \frac{1}{V} \cdot \mathbb{I}(x \in B_i) \frac{N_i}{N} \\ &= \begin{cases} \frac{1}{V} \cdot \frac{N_i}{N}, & x \in B_i \\ \frac{1}{V} \cdot \frac{N_k}{N}, & x \in B_k \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$



$k=3, V=1$ for each bin



$k=2, V=1.5$ for each bin

Note: larger k , smaller bin ($V \downarrow$) \Rightarrow More granularity \Rightarrow overfit

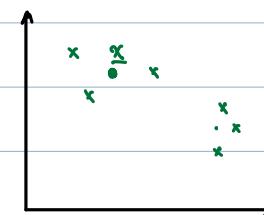
Lecture 2

2) Nearest Neighbor Estimation

For $\forall x \in \mathbb{R}^d$ (not necessary a sample point),

let's define $x_0, x_1, \dots, x_{[k]}$ be the first k

nearest neighbor of x , sorted in ascending order.



let $r_k(x) = \|x - x_{[k]}\|$ (radius of smallest sphere)

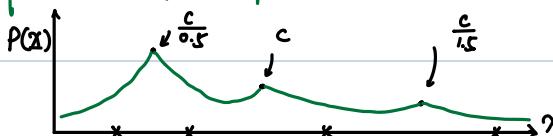
$V_k(x) = \text{volume of sphere in } \mathbb{R}^d \text{ with radius } r_k(x)$

then $\hat{P}(x) \propto \frac{1}{V_k(x)}$

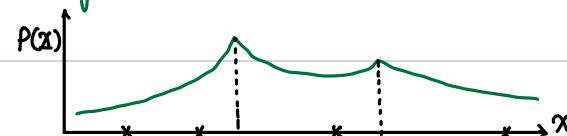
$$\hat{P}(x) = \frac{c}{V_k(x)}$$

To find c , we have $\int \hat{P}(x) = 1$

Example: $D = \{1, 2, 4, 7\}$. set $k=2$.

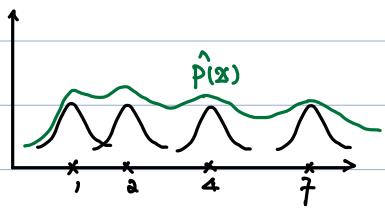


if set $k=3$.



Observation: larger K , smoother output.

3) Parzen Window Estimation with Gaussian kernels



- Each sample has a Gaussian Distribution.
- Add all the probability
- Scale to have CDF = 1

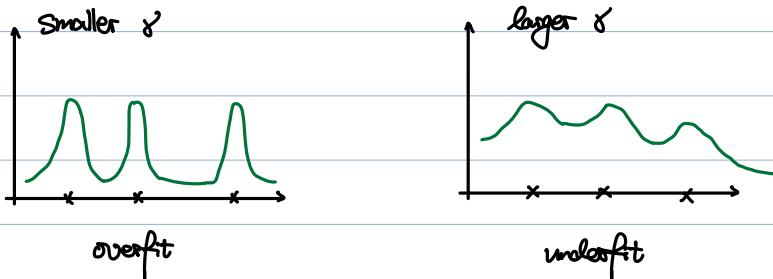
Details: Standard normal distribution (PDF)

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

$$\hat{P}(x) = \frac{1}{K} \sum_{i=1}^N \phi\left(\frac{\|x - x_i\|}{\gamma}\right)$$

γ : kernel width, width of each sample affected neighbor.

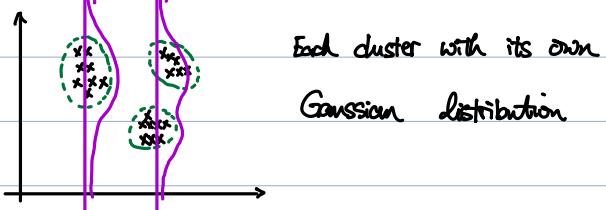
K : Normalizing constant



Note: If more points → complex computation.

Lec 22.

4) Gaussian Mixture Model



It suffices to use only k Gaussian distributions, one fit for each cluster.

- mean: $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^d$

covariance: $\Sigma_1, \Sigma_2, \dots, \Sigma_k \in \mathbb{R}^{d \times d}$

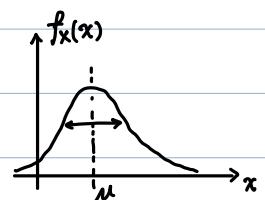
Review:

$d=1$: PDF of Gaussian R.V. X with mean μ , and variance σ^2 .

$$f_X(x) \equiv N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

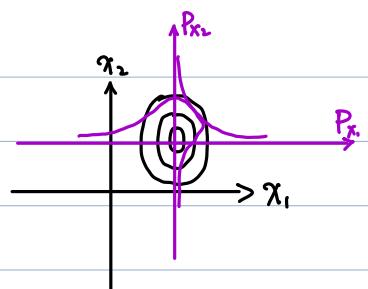
$d=2$: PDF of jointed Gaussian R.V. X_1 and X_2 , with mean μ_1 & μ_2 ,

variances σ_1^2 & σ_2^2 , and correlation coefficient ρ .



$$P = \frac{E[(X_1 - \mu_1)(X_2 - \mu_2)]}{\sigma_1 \cdot \sigma_2}$$

$$f_{X_1, X_2}(x_1, x_2) = N(x_1; \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) \\ = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \cdot e^{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x_1 - \mu_1}{\sigma_1}\right)^2 - \frac{2\rho(x_1 - \mu_1)(x_2 - \mu_2)}{\sigma_1\sigma_2} + \left(\frac{x_2 - \mu_2}{\sigma_2}\right)^2\right]}$$



General d>1:

Consider d jointly Gaussian RVs, $X_1, X_2 \dots X_d$, with $\mu_1, \mu_2 \dots \mu_d$ and covariance matrix.

$$\Sigma \triangleq \begin{bmatrix} \text{cov}(X_1, X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_d) \\ \vdots & \ddots & \ddots & \vdots \\ \text{cov}(X_d, X_1) & \dots & \dots & \text{cov}(X_d, X_d) \end{bmatrix}$$

$$\text{where } \text{cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)]$$

$$\Leftrightarrow \Sigma = E[(X - \mu)(X - \mu)^T]$$

$$f_{\underline{X}}(\underline{x}) = N(\underline{x}, \mu, \Sigma) \\ = \frac{1}{(2\pi)^{\frac{d}{2}} \cdot \det(\Sigma)^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(\underline{x} - \mu)^T \Sigma^{-1} (\underline{x} - \mu)}$$

Data Generalization Model (Hypothesis)

(for each data point in D)

First, pick a bump (cluster)

$j \in \{1, 2, \dots, k\}$ according to probability distribution $\{w_1, w_2 \dots, w_k\}$

$$\begin{cases} w_i > 0 \\ \sum_{i=1}^k w_i = 1 \end{cases}$$

Then, generate \underline{x} according to probability density $P(\underline{x}|j) = N(\underline{x}; \mu_j, \Sigma_j)$

$$= \frac{1}{(2\pi)^{\frac{d}{2}} \cdot \det(\Sigma_j)^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(\underline{x} - \mu_j)^T \Sigma_j^{-1} (\underline{x} - \mu_j)}$$

\Rightarrow the overall PDF for \underline{x} is

$$P(\underline{x}) = \sum_{j=1}^k w_j P(\underline{x}|j) = \sum_{j=1}^k w_j \cdot N(\underline{x}; \mu_j, \Sigma_j)$$

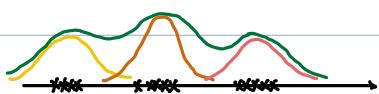
Example: $d=1, k=3$

$$D = \{\underline{x}_1, \underline{x}_2 \dots, \underline{x}_N\}$$

For each index n,



Given D, need to estimate μ_1, μ_2, μ_3 & $\sigma_1, \sigma_2, \sigma_3$.



Given $D = \{\underline{x}_1, \dots, \underline{x}_N\}$, and k . find a GMM $\Omega = \{w_j, \mu_j, \Sigma_j\}_{j=1}^k$ that is the best fit for D.

Best-fit: MLE

$$\hat{P}_\Omega(D) = \prod_{n=1}^N \hat{P}(\underline{x}_n)$$

is maximized.

$$\text{Minimize } E_\Omega(\Omega) \triangleq -\log \hat{P}_\Omega(D)$$

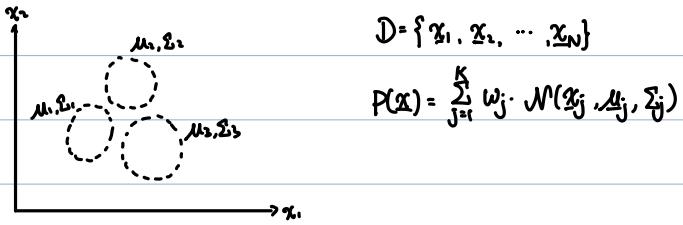
$$= -\sum_{n=1}^N \log \hat{P}(\underline{x}_n)$$

$$E_n(\Omega) = -\log \hat{P}(Z_n)$$

How to minimize $E_n(\Omega)$?

SGD?

Lecture 23



Given set D and k , find the GMM $\Omega = \{w_j, \mu_j, \Sigma_j\}_{j=1}^k$ that is the "best fit" for D .

Maximize the likelihood

$$\hat{P}_\Omega(D) = \prod_{n=1}^N \hat{P}(x_n)$$

$$\equiv \text{Minimize } E_n(\Omega) = -\log \hat{P}_\Omega(D)$$

$$= \sum_{n=1}^N -\log \hat{P}(x_n)$$

$$= \sum_{n=1}^N -\log \underbrace{\sum_{j=1}^k w_j \cdot N(x_n; \mu_j, \Sigma_j)}_{E_n(\Omega)}$$

Can use SGD? No

Complicated E_n

$\sum w_j = 1$. constrained optimization

The Expectation Maximization (EM) Algorithm

An alternating optimization approach with 2 subproblems.

Subproblem 1: (M step)

Given $D = \{x_1, \dots, x_N\}$, suppose for each $x_i \in D$, we know the bump (Gaussian distribution) it belongs to.

Let $B_j \in D$ denotes the j th bump.

i.e. all points in B_j are sampled from

$$P(x|j) = N(x; \mu_j, \Sigma_j)$$

Estimate w_j, μ_j, Σ_j for $j=1, 2, \dots, k$

$$w_j = \frac{N_j}{N} . N_j: \# \text{of points in } B_j$$

$$\mu_j = \frac{1}{N_j} \sum_{x_n \in B_j} x_n \quad (\text{sample mean})$$

$$\Sigma_j = \frac{1}{N_j} \sum_{x_n \in B_j} (x_n - \mu_j)(x_n - \mu_j)^T \quad (\text{sample covariance})$$

Subproblem 2: (E step)

Given $\Omega = \{w_j, \mu_j, \Sigma_j\}_{j=1}^k$, estimate bump membership. i.e. for each \underline{x}_n , find bump j that is most likely to produce \underline{x}_n .

$$\underline{x}_n \in B_{j*} \text{ if } j^* = \operatorname{argmax}_{j \in \{1, 2, \dots, k\}} P_r(j | \underline{x}_n)$$

Recall: Bayes' Thm

$$P_r(j | \underline{x}_n) = \frac{P_r(\underline{x}_n | j) \cdot P_r(j)}{P_r(\underline{x}_n)} = \frac{P_r(\underline{x}_n | j) \cdot P_r(j)}{\sum_{i=1}^k P_r(\underline{x}_n | i) \cdot P_r(i)} = \frac{N(\underline{x}_n; \mu_j, \Sigma_j) \cdot w_j}{\sum_{i=1}^k N(\underline{x}_n; \mu_i, \Sigma_i) \cdot w_i}$$

$$\Rightarrow j^* = \operatorname{argmax}_j N(\underline{x}_n; \mu_j, \Sigma_j) \cdot w_j$$

EM algo. Summary: (Hard decisions)

1. Initialization: Start with arbitrary bump membership for each \underline{x}_n .
2. Estimate $\Omega = \{w_j, \mu_j, \Sigma_j\}_{j=1}^k$ given bump membership $\{B_1, \dots, B_k\}$ (subproblem 1)
3. Estimate bump membership $\{B_1, \dots, B_k\}$ given Ω (subproblem 2)

Repeat step 2 & 3 until convergence.

Note: ① Convergence is guaranteed.

② B_1, \dots, B_k are auxiliary ("hidden" variables).