

CS 480



Personal Notes

Marcus Chan

Taught by Hongyang Zhang

uw cs '25



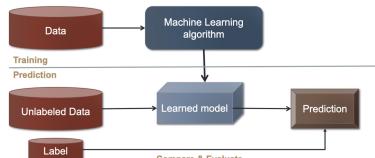
Chapter 1: Perceptrons

ML

Q1 "Machine learning" is a branch of AI that focuses on methods that learn from data & make predictions on unseen data.

Q2 3 phases:

- ① training;
- ② prediction; &
- ③ evaluation.



PARADIGMS OF ML ALGOS (TRAINING)

Q1 "Supervised model": learning with labelled data (x, y)
eg email classification, image classification

Q2 "Unsupervised model": discover patterns in unlabeled data x
eg cluster similar data points, reduce data dimension
etc

Q3 "Semi-supervised model": using both labelled & unlabelled data

WHAT A DATASET LOOKS LIKE

	Training samples					Test samples		
	x_1	x_2	x_3	x_4	\dots	x_n	x'_1	x'_2
$\mathbb{R}^d \ni \text{Feature}$	0	1	0	1	\dots	1	1	0.9
	0	0	1	1	\dots	0	1	1.1
	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\vdots
	1	0	1	0	\dots	1	1	-0.1
Label y	+	+	-	+	\dots	-	?	?

- each column is a data point, n in total & each with d features
- y is the "label vector"
- x' & x'_2 are the test samples whose labels need to be predicted.
- (we use " x " to denote test samples)

INNER PRODUCT: $\langle x, w \rangle$

Q1 Define the "inner product" of a & b to be

$$\langle a, b \rangle = \sum_j a_j b_j,$$

where a_j, b_j are the j^{th} entries of a & b .

LINEAR FUNCTION

Q1 We say a function f is "linear" if

$$f(\alpha x + \beta z) = \alpha f(x) + \beta f(z) \quad \forall \alpha, \beta \in \mathbb{R}, x, z \in \mathbb{R}^d.$$

Q2 Equivalently, f is linear iff there exists $w \in \mathbb{R}^d$ such that

$$f(x) = \langle x, w \rangle = \sum_j x_j w_j.$$

Proof: (\Rightarrow) Let $w = [f(e_1), \dots, f(e_d)]$, where e_i is the i^{th} coordinate vector. Then

$$\begin{aligned} f(x) &= f(x_1 e_1 + \dots + x_d e_d) \\ &= x_1 f(e_1) + \dots + x_d f(e_d) \\ &= \langle x, w \rangle. \end{aligned}$$

(\Leftarrow) Note

$$\begin{aligned} f(\alpha x + \beta z) &= \langle (\alpha x + \beta z), w \rangle \\ &= \alpha \langle x, w \rangle + \beta \langle z, w \rangle \\ &= \alpha f(x) + \beta f(z). \quad \square \end{aligned}$$

AFFINE FUNCTION

Q1 We say f is an "affine function" if there exists a $w \in \mathbb{R}^d$, $b \in \mathbb{R}$ such that

$$f(x) = \langle x, w \rangle + b \quad \forall x \in \mathbb{R}^d.$$

SCORE: \hat{y}

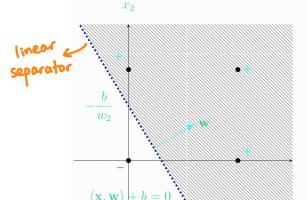
Q1 Given $w \in \mathbb{R}^d$, $b \in \mathbb{R}$, define the "score" at some $x \in \mathbb{R}^d$ to be

$$\text{Score}(x) = \langle x, w \rangle + b.$$

Q2 Our "prediction" for y is then

$$\hat{y} = \text{sign}(\text{Score}(x)) = \begin{cases} +1, & \text{Score}(x) > 0 \\ -1, & \text{Score}(x) \leq 0. \end{cases}$$

We want to tune w, b so that " $\hat{y} = y$ " for each x .



- x is free, w & b fixed
- w & b uniquely determine the linear separator.

PERCEPTRONS

Algorithm for training:

Algorithm 1 Training Perceptron

Input: Dataset = $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \{\pm 1\}$: $i = 1, \dots, n$, initialization $\mathbf{w}_0 \in \mathbb{R}^d$ and $b_0 \in \mathbb{R}$

Output: \mathbf{w} and b (so a linear classifier $\text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle + b)$)

for $t = 1, 2, \dots$ do

- receive index $I_t \in \{1, \dots, n\}$ // I_t can be random
- if $y_{I_t}(\langle \mathbf{x}_{I_t}, \mathbf{w} \rangle + b) \leq 0$ // a "mistake" happens
- then

 - $\mathbf{w} \leftarrow \mathbf{w} + y_{I_t} \mathbf{x}_{I_t}$ // update after a "mistake"
 - $b \leftarrow b + y_{I_t}$

- end

end

- we typically set $w_0=0$ & $b_0=0$

- we only update after a mistake

(aka "lazy update")

- note we are going through the data one by one.

Q2 In particular, we want to find $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$ such that for all $i=1, \dots, n$,

$$y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) > 0.$$

Q3 Note that if a mistake happens on (x, y) :

$$\begin{aligned} y[\langle \mathbf{x}, \mathbf{w}_{k+1} \rangle + b_{k+1}] &= y[\langle \mathbf{x}, \mathbf{w}_k + y \mathbf{x} \rangle + b_k + y] \\ &= y[\langle \mathbf{x}, \mathbf{w}_k \rangle + y \langle \mathbf{x}, \mathbf{x} \rangle + b_k + y] \\ &= y[\langle \mathbf{x}, \mathbf{w}_k \rangle + y \|\mathbf{x}\|_2^2 + b_k + y] \\ &= y[\langle \mathbf{x}, \mathbf{w}_k \rangle + b_k] + y^2 \|\mathbf{x}\|_2^2 + y^2 \\ &= y[\langle \mathbf{x}, \mathbf{w}_k \rangle + b_k] + \underbrace{\|\mathbf{x}\|_2^2 + 1}_{\text{always positive } \& \geq 1}. \end{aligned}$$

Q4 Example: spam filtering.

	x_1	x_2	x_3	x_4	x_5	x_6
and	1	0	0	1	1	1
viagra	1	0	1	0	0	0
the	0	1	1	0	1	1
of	1	1	0	1	0	1
nigeria	1	0	0	0	1	0
y	+	-	+	-	+	-

- Recall the update: $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$, $b \leftarrow b + y$ (when a mistake happens on (x, y))

- $\mathbf{w}_0 = [0, 0, 0, 0, 0]$, $b_0 = 0 \implies \text{score}(\mathbf{x}_1) = 0 \implies \hat{y}_1 = - \quad \times$
- $\mathbf{w}_1 = [1, 1, 0, 1, 1]$, $b_1 = 1 \implies \text{score}(\mathbf{x}_2) = 2 \implies \hat{y}_2 = + \quad \times$
- $\mathbf{w}_2 = [1, 1, -1, 0, 1]$, $b_2 = 0 \implies \text{score}(\mathbf{x}_3) = 0 \implies \hat{y}_3 = - \quad \times$
- $\mathbf{w}_3 = [1, 2, 0, 0, 1]$, $b_3 = 1 \implies \text{score}(\mathbf{x}_4) = 2 \implies \hat{y}_4 = + \quad \times$
- $\mathbf{w}_4 = [0, 2, 0, -1, 1]$, $b_4 = 0 \implies \text{score}(\mathbf{x}_5) = 1 \implies \hat{y}_5 = + \quad \checkmark$
- $\mathbf{w}_4 = [0, 2, 0, -1, 1]$, $b_4 = 0 \implies \text{score}(\mathbf{x}_6) = -1 \implies \hat{y}_6 = - \quad \checkmark$

13 / 2

A TRICK TO HIDE THE BIAS TERM

Q1 Note that

$$\langle \mathbf{x}, \mathbf{w} \rangle + b = \underbrace{\langle (\mathbf{x})_1, (\mathbf{w})_1 \rangle}_{\mathbf{x}_{\text{pad}}} + \underbrace{\langle (\mathbf{x})_2, (\mathbf{w})_2 \rangle}_{\mathbf{w}_{\text{pad}}}$$

This is a "trick" to ignore b in future calculations.

Q2 Thus, our new update rule is

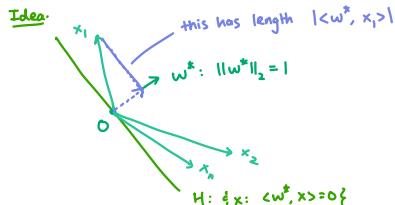
$$\mathbf{w}_{\text{pad}} \leftarrow \mathbf{w}_{\text{pad}} + y \mathbf{x}_{\text{pad}}.$$

CONVERGENCE THEOREM (LINEARLY SEPARABLE CASE)

\bullet Suppose there exists a w^* such that
 $y_i \langle x_i, w^* \rangle > 0 \quad \forall i=1, \dots, n.$

Assume $\|x_i\|_2 \leq C$ and that w^* is normalized so that $\|w^*\|_2 = 1$.

Define the margin $\gamma = \min_i |\langle x_i, w^* \rangle|$. Then the Perceptron algorithm converges after C^2/γ^2 mistakes.



- w^* is our "perfect" solution for w (ie the "goal" criteria is satisfied).
- thus, we want to show w "converges" to w^* .

Proof. Recall the update is $w \leftarrow w + yx$.

Define

$$\cos(w, w^*) = \frac{\langle w, w^* \rangle}{\|w\| \|w^*\|} = \frac{\langle w, w^* \rangle}{\|w\|} \quad (\text{since we defined } \|w^*\|=1).$$

Consider an update and its effect on $\langle w, w^* \rangle$:

$$\begin{aligned} \langle w, w^* \rangle &\longrightarrow \langle w+yx, w^* \rangle \\ &= \langle w, w^* \rangle + y \underbrace{\langle x, w^* \rangle}_{\text{positive } \because w^* \text{ is perfect}} \\ &= \langle w, w^* \rangle + |\langle x, w^* \rangle| \\ &\geq \langle w, w^* \rangle + \gamma. \end{aligned}$$

This means for each update, $\langle w, w^* \rangle$ grows by at least $\gamma > 0$.

Similarly, consider an update's effect on $\|w\|_2^2$:

$$\begin{aligned} \|w\|_2^2 &= \langle w, w \rangle \longrightarrow \langle w+yx, w+yx \rangle \\ &= \langle w, w \rangle + 2y \underbrace{\langle x, x \rangle}_{<0} + \\ &\quad y^2 \langle x, x \rangle \\ &= \langle w, w \rangle + 2y \langle w, x \rangle + \underbrace{\|x\|_2^2}_{\leq C^2} \\ &\leq \langle w, w \rangle + C^2. \end{aligned}$$

This means for each update, $\langle w, w \rangle$ grows by at most C^2 .

Now, let $w_0 = 0$. We now know after M updates:

$$\begin{aligned} \langle w_M, w^* \rangle &\geq \langle w_{M-1}, w^* \rangle + \gamma \\ &\geq \langle w_{M-2}, w^* \rangle + 2\gamma \\ &\geq \dots \geq \underbrace{\langle w_0, w^* \rangle}_{=0} + M\gamma \\ &= M\gamma. \end{aligned}$$

Similarly, note

$$\begin{aligned} \langle w_M, w_M \rangle &\leq \langle w_{M-1}, w_{M-1} \rangle + C^2 \\ &\leq \dots \leq \underbrace{\langle w_0, w_0 \rangle}_{=0} + MC^2 \\ &\leq MC^2. \end{aligned}$$

Since

$$\cos(w, w^*) = \frac{\langle w, w^* \rangle}{\|w\|} \leq 1 \Rightarrow \langle w, w^* \rangle \leq \|w\|$$

Therefore

$$M\gamma \leq \langle w, w^* \rangle \leq \|w\| \leq \sqrt{MC^2} = \sqrt{M}C.$$

Rearranging, this tells us that $M \leq \frac{C^2}{\gamma^2}$, which finishes the proof. \square

\bullet In particular, the larger γ is, the more separable the data is, and hence the faster the algorithm converges!

ANOTHER PERSPECTIVE ON PERCEPTRONS

Our hypothesis is $\hat{y} = \text{sign}\{\langle w, x \rangle\}$.

We can define our "loss function" as

$$\begin{aligned} l(w; x_t; y_t) &= -y_t \langle w, x_t \rangle \mathbb{I}[\text{mistake on } (x_t, y_t)] \\ &= \begin{cases} -y_t \langle w, x_t \rangle, & \text{if mistake happens} \\ 0 & \Leftrightarrow y_t \langle w, x_t \rangle < 0 \\ 0, & \text{otherwise} \end{cases} \\ &= -\min\{0, y_t \langle w, x_t \rangle\}. \end{aligned}$$

The average of all the loss functions of the data points is then

$$L(w) = -\frac{1}{n} \sum_{t=1}^n y_t \langle w, x_t \rangle \mathbb{I}[\text{mistake on } x_t].$$

Our gradient descent update:

$$w_{t+1} = w_t - \gamma_t \nabla_w l(w_t, x_t, y_t) = w_t + \gamma_t y_t x_t \mathbb{I}[\text{mistake on } x_t].$$

If we set the step size $\gamma_t = 1$, then

$$w_{t+1} = w_t + y_t x_t,$$

which is our update rule.

PERCEPTRONS ARE NOT UNIQUE

Note perceptrons are not unique as the algorithm terminates as long as there is no mistake.

- it depends on initialization & our sampling rule of I_t .

MAXIMIZE MARGIN

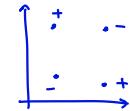
We want to choose w such that

$$w = \max_{\substack{w: \forall i, \hat{y}_i y_i > 0}} \min_{i=1, \dots, n} \frac{\hat{y}_i y_i}{\|w\|}, \quad \hat{y}_i := \langle x_i, w \rangle + b.$$

XOR DATASET

There is no line that can separate + from -.

x ₁	x ₂	x ₃	x ₄
0	1	0	1
0	0	1	1
-	+	+	-



What if we run Perceptron?

Suppose $\exists w, b$ s.t. $y(\langle x, w \rangle + b) > 0$. Then:

$$x_1 = (0,0), y_1 = - \Rightarrow b < 0$$

$$x_2 = (1,0), y_2 = + \Rightarrow w_1 + b > 0 \quad \left\{ \begin{array}{l} w_1 + w_2 + 2b > 0 \\ w_2 + b > 0 \end{array} \right. \quad > 0$$

$$x_3 = (0,1), y_3 = + \Rightarrow w_2 + b > 0 \quad > 0$$

$$x_4 = (1,1), y_4 = - \Rightarrow w_1 + w_2 + 2b < 0. \quad < 0$$

Hence

$$\underbrace{(w_1 + w_2 + 2b)}_{> 0} - \underbrace{(w_1 + w_2 + b)}_{< 0} = b > 0,$$

which contradicts our earlier statement that $b < 0$.

HARDNESS RESULT (NON-LINEARLY SEPARABLE CASE)

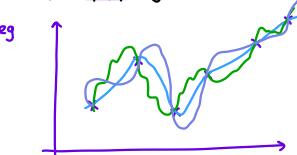
If there is no perfect separating hyperplane for our data, then the Perceptron algorithm cycles.

Chapter 2: Linear Regression

Q₁ Idea: Given training data (x_i, y_i) , find a $f: X \rightarrow Y$ such that $f(x_i) \approx y_i$, where

- ① $x_i \in X \subseteq \mathbb{R}^d$: the feature vector for the i^{th} training example
- ② $y_i \in Y \subseteq \mathbb{R}^t$: t responses
 - note we could have $t=1$ or even $t=\infty$

Q₂ Note for any finite training data (x_i, y_i) , $i=1, \dots, n$, there exist infinitely many functions f such that for all i , $f(x_i) = y_i$.



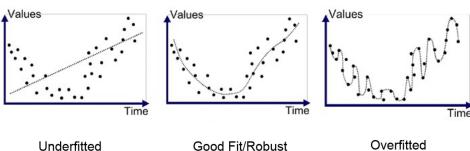
Q₃ Moreover, our prediction $\hat{y} = f(x)$ can vary significantly on new data x !

Q₄ To choose f , we can

- ① leverage prior knowledge of f ; & eg if x & y come from a population which follows "rules"
- ② choose the "simplest" function.

UNDERFITTING, GOOD FITTING,

OVERFITTING



STATISTICAL LEARNING

Q We assume the training & test data are both iid samples from the same unknown distribution P ; ie

$$(X_i, Y_i) \sim P$$
$$(X, Y) \sim P.$$

LEAST SQUARES REGRESSION

Q We want to choose f so that

$$f = \min_{f: X \rightarrow Y} E \|f(x) - y\|_2^2.$$

this is our least squared error.

REGRESSION FUNCTION: $m(x)$

Our "regression function" is

$$f^*(x) = m(x) = \mathbb{E}[y | x=x].$$

However, calculating m requires us to know the distribution of P , ie all pairs (X, Y) .

We show that m is optimal; ie

$$m(x) = \min_{f: \mathbb{R} \rightarrow \mathbb{R}} \mathbb{E}_{x \sim P} \|f(x) - y\|_2^2.$$

Proof. First, see that

$$\begin{aligned} \mathbb{E} \|f(x) - y\|_2^2 &= \mathbb{E} \|\mathbb{E}[f(x) - m(x) + m(x) - y]\|_2^2 \\ &= \mathbb{E} \|f(x) - m(x)\|_2^2 + \mathbb{E} \|m(x) - y\|_2^2 \\ &\quad + 2 \mathbb{E} \langle f(x) - m(x), m(x) - y \rangle. \end{aligned}$$

Using $\|ab\|_2^2 = \|a\|_2^2 + \|b\|_2^2 + 2\langle a, b \rangle$

Then

$$\begin{aligned} \mathbb{E}_{x,y} [\langle f(x) - m(x), m(x) - y \rangle] &= \mathbb{E}_x [\mathbb{E}_{y|x} [\langle f(x) - m(x), m(x) - y \rangle]] \\ &\quad (\text{by double expectation theorem, see STAT 330}) \\ &= \mathbb{E}_x [\langle f(x) - m(x), m(x) - \mathbb{E}[y|x] \rangle] \\ &= \mathbb{E}_x [\langle f(x) - m(x), 0 \rangle] \\ &= 0. \end{aligned}$$

Hence

$$\mathbb{E} \|f(x) - y\|_2^2 = \mathbb{E} \|f(x) - m(x)\|_2^2 + \underbrace{\mathbb{E} \|m(x) - y\|_2^2}_{\text{noise (variance) term}}.$$

- independent wrt f .

Therefore, to reduce $\mathbb{E} \|f(x) - y\|_2^2$, we need to only minimize $\mathbb{E} \|f(x) - m(x)\|_2^2$, which is minimal (ie = 0) when $f = m$!

However, m is unaccessible since the conditional distribution is unknown, so we need to try to get close to m using the training data.

BIAS-VARIANCE TRADEOFF

Let f_D be the regressor learned on the training dataset D . Then

$$\begin{aligned} \mathbb{E}_{D, X, Y} \|f_D(x) - y\|_2^2 &= \mathbb{E}_x \|\mathbb{E}_D [f_D(x)] - m(x)\|_2^2 \\ &\quad \underbrace{\text{test error}}_{\text{bias}^2} \\ &\quad + \mathbb{E}_{D, X} \|f_D(x) - \mathbb{E}_D [f_D(x)]\|_2^2 \\ &\quad \underbrace{\text{variance}}_{\text{noise}} \\ &\quad + \mathbb{E}_{X, Y} \|m(x) - y\|_2^2 \end{aligned}$$

Proof. We have shown

$$\begin{aligned} \mathbb{E}_{X, Y} \|f_D(x) - y\|_2^2 &= \mathbb{E}_x \|\mathbb{E}_D [f_D(x)] - m(x)\|_2^2 \\ &\quad + \underbrace{\mathbb{E}_{X, Y} \|m(x) - y\|_2^2}_{\text{noise - independent wrt } f_D}. \end{aligned}$$

Taking E_D of both sides:

$$\begin{aligned} \mathbb{E}_D \mathbb{E}_{X, Y} \|f_D(x) - y\|_2^2 &= \mathbb{E}_D \mathbb{E}_X \|\mathbb{E}_D [f_D(x)] - m(x)\|_2^2 \\ &\quad + \mathbb{E}_{X, Y} \|m(x) - y\|_2^2. \quad \text{①} \end{aligned}$$

Define $\bar{f}(x) = \mathbb{E}_D [f_D(x)]$.

Idea: We can sample multiple f 's from various samples D :

$$\begin{aligned} D_i \sim P &\rightarrow f_{D_i} \\ &\vdots \\ D_n \sim P &\rightarrow f_{D_n} \end{aligned} \quad \left\{ \begin{array}{l} \text{then we define} \\ \bar{f}(x) = \text{avg } f_{D_i}(x). \end{array} \right.$$

Then

$$\begin{aligned} \mathbb{E}_D \mathbb{E}_X \|f_D(x) - m(x)\|_2^2 &= \mathbb{E}_{D, X} \|f_D(x) - \bar{f}(x) + \bar{f}(x) - m(x)\|_2^2 \\ &= \mathbb{E}_{D, X} \|f_D(x) - \bar{f}(x)\|_2^2 + \mathbb{E}_{D, X} \|\bar{f}(x) - m(x)\|_2^2 \\ &\quad + 2 \mathbb{E}_{D, X} \langle f_D(x) - \bar{f}(x), \bar{f}(x) - m(x) \rangle. \end{aligned}$$

Similarly, see that

$$\begin{aligned} \mathbb{E}_{D, X} \langle \bar{f}(x) - f_D(x), m(x) - \bar{f}(x) \rangle &= \mathbb{E}_X \mathbb{E}_D \langle m(x) - \bar{f}(x), \bar{f}(x) - f_D(x) \rangle \\ &\quad \text{constant wrt } D \\ &= \mathbb{E}_X \langle m(x) - \bar{f}(x), \bar{f}(x) - \underbrace{\mathbb{E}_D [f_D(x)]}_{\bar{f}(x)} \rangle \\ &= 0. \end{aligned}$$

Expanding ① yields the result desired. \blacksquare

In particular, as the model capacity increases,

- ① the bias term decreases (ie model is more expressively powerful); but
- ② the variance increases (ie model is less stable).

SAMPLING → TRAINING

In practice, we can only calculate the sample average, ie we find f so that

$$f = \min_{f: X \rightarrow Y} \hat{E} \|f(X) - Y\|_2^2 := \frac{1}{n} \sum_{i=1}^n \|f(x_i) - y_i\|_2^2.$$

However, as our training data size $n \rightarrow \infty$, $\hat{E} \rightarrow E$ & hopefully $\operatorname{argmin} \hat{E} \rightarrow \operatorname{argmin} E$.

LINEAR REGRESSION

In linear regression, our regression functions are "affine"; ie in the form

$$f(x) = Wx + b, \quad W \in \mathbb{R}^{t \times d}, \quad b \in \mathbb{R}^t.$$

- $t = \#$ of response parameters we want to predict
- $d = \#$ of input parameters

Again, we can use padding:

$$x \leftarrow \begin{pmatrix} x \\ 1 \end{pmatrix}, \quad w \in [w, b] \Rightarrow f(x) = Wx$$

In matrix form:

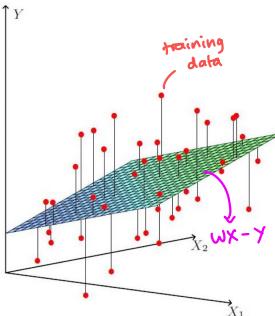
$$\frac{1}{n} \sum_i \|f(x_i) - y_i\|_2^2 = \frac{1}{n} \|Wx - y\|_F^2,$$

$$X \in [\dots, x_n] \in \mathbb{R}^{(d+1) \times n}, \quad Y = [y_1, \dots, y_n] \in \mathbb{R}^{t \times n},$$

$$\|A\|_F = \sqrt{\sum_{i,j} a_{ij}^2}$$

We want to find W such that

$$W = \min_{W \in \mathbb{R}^{t \times (d+1)}} \frac{1}{n} \|Wx - y\|_F^2.$$



- geometrically, we want to minimise the sum of distances between the input training data & the resultant hyperplane.

SOLVING LINEAR REGRESSION

We define our loss function as

$$\text{Loss}(W) = \frac{1}{n} \|Wx - y\|_F^2$$

Taking the derivative wrt W & setting to zero:

$$\begin{aligned} \nabla_W \text{Loss}(W) &= \frac{2}{n} (Wx - y) x^T (= 0) \\ \Rightarrow Wx x^T &= y x^T \\ \Rightarrow W &= y x^T (x x^T)^{-1} \end{aligned}$$

PREDICTION

Once we have solved W on the training set (X, Y) , we can predict on unseen data x_{test} :

$$\hat{Y}_{\text{test}} = W x_{\text{test}}$$

The "test error" (if true labels were available) is

$$\text{test error} = \frac{1}{n_{\text{test}}} \|y_{\text{test}} - \hat{Y}_{\text{test}}\|_F^2$$

The "training error" is

$$\text{training error} = \frac{1}{n} \|y - Wx\|_F^2.$$

We can minimize the training error to reduce the test error.

ILL-CONDITIONING

Consider $X = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$, $y = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$. Solving linear least squares regression:

$$w = y x^T (x x^T)^{-1} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \begin{pmatrix} 1/3 & 1/3 \\ 1/3 & 1 \end{pmatrix} = \begin{pmatrix} -2/3 \\ 1 \end{pmatrix}$$

So slight perturbation leads to chaotic behavior!

This occurs when X is ill-conditioned; ie close to rank deficient.

- two cols in X are close to linearly dependent
- but corresponding y 's are different
- this is a contradiction $\Rightarrow w$ becomes unstable.

RIDGE REGRESSION

Idea: We instead try to find

$$W = \min_W \left[\frac{1}{n} \|WX - Y\|_F^2 + \lambda \|W\|_F^2 \right]$$

Why is this better?

consider Loss(W) = $\frac{1}{n} \|WX - Y\|_F^2 + \lambda \|W\|_F^2$.

$$\Rightarrow \nabla_W \text{Loss}(W) = \frac{2}{n} (WX - Y) X^T + 2\lambda W (= 0)$$

$$WX X^T - Y X^T + 2\lambda W = 0$$

$$WX X^T - Y X^T + W(2\lambda I) = 0$$

$$WX X^T + W(2\lambda I) = Y X^T$$

$$\therefore W = (X X^T + 2\lambda I)^{-1} (Y X^T)$$

Then $XX^T + n\lambda I$ is far from rank-deficient matrices for large λ . (Proof uses SVD - see MATH 235).

② controls our trade-off:

① $\lambda=0$ reduces to ordinary linear regression;

② $\lambda=\infty$ reduces to $W=0$;

③ intermediate λ restricts output to be

$\frac{1}{\lambda}$ proportional to input.

Alternatively, note

$$\frac{1}{n} \|WX - Y\|_F^2 + \lambda \|W\|_F^2 = \frac{1}{n} \|W[X \sqrt{n\lambda I}] - [Y \mathbf{0}]\|_F^2$$

So we can also

① augment X with $\sqrt{n\lambda I}$; ie $\tilde{X} = (X \sqrt{n\lambda I})$

② augment Y with zeroes; ie $\tilde{Y} = (Y \mathbf{0})$

(ie data augmentation) to achieve regularization.

Chapter 3: Logistic Regression

MOTIVATION

Q₁: This is for linear classification.

Q₂: We can use $\langle \mathbf{x}, \mathbf{w} \rangle$ (our margin) as a measure of our confidence in the prediction \hat{y} .

Q₃: However, as this is un-normalized, it is hard to interpret.

MLE

Q₁: We want to directly learn our "confidence".

$$p(\mathbf{x}; \mathbf{w}) := P(Y=1 | \mathbf{X}=\mathbf{x})$$

Q₂: Then, if $y_1, \dots, y_n, x_1, \dots, x_n$ are independent, then

$$\begin{aligned} & P(Y_1=y_1, \dots, Y_n=y_n | X_1=x_1, \dots, X_n=x_n) \\ &= \prod_{i=1}^n P(Y_i=y_i | X_i=x_i) \\ &= \prod_{i=1}^n [p(x_i; \mathbf{w})]^{y_i} [1-p(x_i; \mathbf{w})]^{1-y_i} \quad \text{if } y_i \in \{0, 1\} \end{aligned}$$

Q₃: Maximizing the likelihood:

$$\begin{aligned} & \max_{\mathbf{w}} \prod_{i=1}^n [p(x_i; \mathbf{w})]^{y_i} [1-p(x_i; \mathbf{w})]^{1-y_i} \\ & \Leftrightarrow \min_{\mathbf{w}} \sum_{i=1}^n [-y_i \log p(x_i; \mathbf{w}) - (1-y_i) \log(1-p(x_i; \mathbf{w}))] \end{aligned}$$

Q₄: We thus want to find \mathbf{w} which satisfies the above optimization problem.

THE LOGIT TRANSFORM

Q₁: If we assume the log of odds ratio is linear: ie

$$\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$$

* we can only perform logistic regression if we assume this!

then

$$p(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$$

↳ this is also called the "sigmoid transformation".

Q₂: Plugging this into the earlier optimization problem, we want to find

$$\min_{\mathbf{w}} \sum_{i=1}^n \log [1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + (1-y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle$$

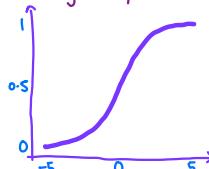
if $y_i \in \{0, 1\}$.

Q₃: If instead $y_i \in \{-1, 1\}$, then

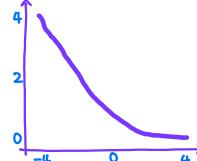
$$\min_{\mathbf{w}} \sum_{i=1}^n \log [1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]$$

↳ this is "logistic loss".

Sigmoid function



logistic loss



TRAINING LOGISTIC REGRESSION

Q Our gradient descent algorithm is

$$w \leftarrow w - \eta \nabla_w \text{Loss}(w)$$

PREDICTION

Q₁ We take

$$\hat{y} = 1 \Leftrightarrow P(Y=1 | X=x) > \frac{1}{2} \Leftrightarrow \langle x, w \rangle > 0$$

Q₂ Our decision boundary is still

$$H := \{x : \langle x, w \rangle = 0\}$$

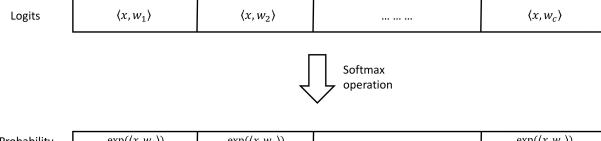
Q₃ So we can predict $\hat{y} = \text{sign}(\langle x, w \rangle)$ as before,
but now with confidence $p(x; w)$.

MULTI-CLASS EXTENSION

Q₁ Idea: For a class $y \in \{1, \dots, c\}$, we want
to learn $\{w_1, \dots, w_c\}$ for each class.

Q₂ We consider the "softmax" function:

$$P(Y=k | X=x, W=[w_1, \dots, w_c]) = \frac{\exp(\langle x, w_k \rangle)}{\sum_{l=1}^c \exp(\langle x, w_l \rangle)}$$



- we map a real-valued vector to a probability vector
- these are non-negative & sum to 1.

Q₃ Training: again, we use MLE:

$$\min_w E \left[-\log \frac{\exp(\langle x, w_y \rangle)}{\sum_{l=1}^c \exp(\langle x, w_l \rangle)} \right]$$

Q₄ Prediction:

$$\hat{y} = \underset{k}{\operatorname{argmax}} P(Y=k | X=x; W=[w_1, \dots, w_c])$$

Chapter 4: Hard-Margin Support Vector Machines

INTRODUCTION

$\hat{y}_i = \hat{w}^T x_i + b$, and don't use padding.

Perceptron: we find any $w \in \mathbb{R}^d$, $b \in \mathbb{R}$ such that

$$\begin{aligned} & \min_{w, b} \text{ s.t. } y_i \hat{y}_i > 0 \quad \forall i, \\ & \hat{y}_i = \hat{w}^T x_i + b \\ & \Leftrightarrow \min_{w, b} \text{ s.t. } y_i \hat{y}_i \geq 1 \quad \forall i \end{aligned}$$

However, the larger the margin, the faster Perceptron converges.

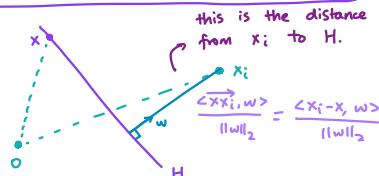
recall # mistakes, $M \leq \frac{C^2}{\gamma^2}$, $\|x_i\|_2 \leq C$, $\gamma = \min_i |\langle x_i, w^* \rangle|$, $\|w^*\|_2 = 1$.

So, the goal of hard-margin SVM is to maximize the margin (assuming data is linearly separable).

DISTANCE FROM A POINT TO A HYPERPLANE

Let $H := \{x : \langle x, w \rangle + b = 0\}$. Then

$$\begin{aligned} \text{distance}(x_i, H) &= \frac{|\langle x_i, w \rangle|}{\|w\|_2}, \quad x \in H \\ &= \frac{|\langle x_i, w \rangle - \langle x, w \rangle|}{\|w\|_2} \\ &= \frac{|\langle x_i, w \rangle + b|}{\|w\|_2} \quad \because x \in H \\ &= \frac{|y_i \hat{y}_i|}{\|w\|_2} \quad \because y_i \hat{y}_i = 0 \end{aligned}$$



MARGIN

We define the "margin" as the smallest distance to a separating hyperplane H among all separable training data; ie

$$\begin{aligned} \text{margin} &= \min_i \frac{|y_i \hat{y}_i|}{\|w\|_2} = \min_i \frac{|\langle x_i, w \rangle + b|}{\|w\|_2}, \\ H &= \{x : \langle x, w \rangle + b = 0\} \end{aligned}$$

eg



Our goal is to maximize the margin among all hyperplanes: ie find

$$\max_{w, b} \min_i \frac{|y_i \hat{y}_i|}{\|w\|_2} \quad \text{s.t. } y_i \hat{y}_i > 0 \quad \forall i$$

TRANSFORMING TO STANDARD FORM

- Q₁: Note for the margin, (w, b) & (cw, cb) has the same loss for $c > 0$.
- Q₂: So, we can fix the numerator arbitrarily to 1:

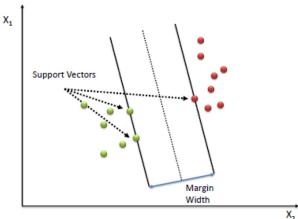
$$\max_{w, b} \left[\frac{1}{\|w\|_2} \text{ s.t. } \min_i y_i \hat{y}_i = 1 \right]$$
$$\Rightarrow \min_{w, b} \left[\frac{1}{2} \|w\|_2^2 \text{ s.t. } y_i (\langle x_i, w \rangle + b) \geq 1 \forall i \right]$$

COMPARISON TO PERCEPTRON

Hard-margin SVM	Perceptron
$\min_{w, b} \frac{1}{2} \ w\ _2^2 \text{ s.t. } y_i \hat{y}_i \geq 1 \forall i$	$\min_{w, b} 0 \text{ s.t. } y_i \hat{y}_i \geq 1 \forall i$
- quadratic programming	- linear programming
- unique solution	- infinitely many solutions
- maximal margin	- convergence rate depends on max margin

SUPPORT VECTORS

- Q₁: Note that
- $$y_i \hat{y}_i \geq 1 \forall i \Leftrightarrow \hat{y}_i \geq +1 \forall i: y_i = +1$$
- $$\hat{y}_i \leq -1 \forall i: y_i = -1$$
- Q₂: This yields 3 parallel hyperplanes:
- $$H = \{x : \langle x, w \rangle + b = 0\}$$
- $$H^+ = \{x : \langle x, w \rangle + b = +1\}$$
- $$H^- = \{x : \langle x, w \rangle + b = -1\}$$
- Q₃: "Support vectors" are those where points lie on the supporting hyperplanes.



LAGRANGIAN DUAL

First, we show

$$\begin{aligned} & \min_{w,b} \frac{1}{2} \|w\|_2^2 \quad \text{s.t. } y_i(\langle x_i, w \rangle + b) \geq 1 \quad \forall i \\ &= \min_{w,b} \max_{\alpha \geq 0} \frac{1}{2} \|w\|_2^2 - \sum_i \alpha_i [y_i(\langle x_i, w \rangle + b) - 1] \\ & \quad \downarrow \\ & \alpha = [\alpha_1, \dots, \alpha_n] \in \mathbb{R}^n; \\ & \alpha \geq 0 \Leftrightarrow \alpha_i \geq 0 \quad \forall i \end{aligned}$$

Proof. let Δ be the second expression.

See that

$$\Delta = \min_{w,b} \max_{\alpha \geq 0} \frac{1}{2} \|w\|_2^2 - \sum_i \alpha_i [y_i(\langle x_i, w \rangle + b) - 1]$$

If $\exists i$ s.t. $y_i(\langle x_i, w \rangle + b) < 1$, then if we set $\alpha_i = \infty$, it follows that $\Delta = +\infty$, which is the maximal value Δ can take.

Otherwise, ie if $\forall i, y_i(\langle x_i, w \rangle + b) \geq 1$,

then

$$\begin{aligned} \Delta &= \frac{1}{2} \|w\|_2^2 - \sum_i \underbrace{\alpha_i}_{\text{tve}} \underbrace{[y_i(\langle x_i, w \rangle + b) - 1]}_{\text{tve}} \\ &\leq \frac{1}{2} \|w\|_2^2. \end{aligned}$$

If we set $\alpha_i = 0 \quad \forall i$, we get $\Delta = \frac{1}{2} \|w\|_2^2$, which is the max value Δ can take.

Therefore,

$$\begin{aligned} \Delta &= \min_{w,b} \begin{cases} +\infty, & \text{if } \exists i \text{ s.t.} \\ & y_i(\langle x_i, w \rangle + b) < 1 \\ \frac{1}{2} \|w\|_2^2, & \text{otherwise} \end{cases} \\ &= \min_{w,b} \frac{1}{2} \|w\|_2^2 \quad \text{if } y_i(\langle x_i, w \rangle + b) \geq 1 \end{aligned}$$

as needed. \square

We can swap the min & max:

$$\max_{\alpha \geq 0} \min_{w,b} \frac{1}{2} \|w\|_2^2 - \sum_i \alpha_i [y_i(\langle x_i, w \rangle + b) - 1]$$

(because of "strong duality")

Now, suppose we fix α , and consider the inner minimization problem.

Then w, b minimizes the function if

$$\frac{\partial}{\partial w} = \frac{\partial}{\partial b} = 0.$$

$$\text{let Loss}(w, b) = \frac{1}{2} \|w\|_2^2 - \sum_i \alpha_i [y_i(\langle x_i, w \rangle + b) - 1].$$

$$\Rightarrow \frac{\partial}{\partial w} = w - \sum_i \alpha_i y_i x_i (= 0), \quad \frac{\partial}{\partial b} = - \sum_i \alpha_i y_i (= 0)$$

$$\rightarrow w = \sum_i \alpha_i y_i x_i, \quad \sum_i \alpha_i y_i = 0.$$

Finally, we consider the "outer" maximization problem.

Plugging in our value of w above:

$$\begin{aligned} \Rightarrow \text{Loss}(\alpha) &= \frac{1}{2} \left\| \sum_i \alpha_i y_i x_i \right\|_2^2 - \left\langle \sum_i \alpha_i y_i x_i, \sum_i \alpha_i y_i \right\rangle \\ &\quad - b \sum_i \alpha_i y_i + \sum_i \alpha_i \\ &= -\frac{1}{2} \left\| \sum_i \alpha_i y_i x_i \right\|_2^2 + \sum_i \alpha_i \quad \text{s.t. } \sum_i \alpha_i y_i = 0 \end{aligned}$$

Thus, our problem becomes

$$\begin{aligned} & \star \max_{\alpha \geq 0} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \quad \text{s.t. } \sum_i \alpha_i y_i = 0 \\ &= \min_{\alpha \geq 0} - \sum_i \alpha_i + \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \quad \text{s.t. } \sum_i \alpha_i y_i = 0 \end{aligned}$$

WHY USE THE DUAL FORM?

Idea: If data is not linearly separable, we use a non-linear mapping ϕ to map the data.

$$\begin{aligned} \min_{\alpha \geq 0} & - \sum_i \alpha_i + \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \phi(x_i), \phi(x_j) \rangle \\ \text{s.t. } & \sum_i \alpha_i y_i = 0. \end{aligned}$$