

- Regression: ended up with optimization problem:

$$\min_w \frac{1}{m} \sum_{i=1}^m (w^T x^i - y^i)^2 = (\text{square loss function}),$$

convex + differentiable,
there can use GD!

$$E[Y|X] = w^T x \quad (*)$$

even if $*$ does not hold, can still find min cost(?).

- Classification:

what is the loss function for classification?

our guess for label $y^i \in \{0, 1\}$ is going to be $\text{sign}(w^T x)$ for some vector w .

- Penalized if $\text{sign}(w^T x^i) \neq y^i$

- No penalty if $\text{sign}(w^T x^i) = y^i$.

Examine quantity $y^i \cdot (w^T x^i)$:

if negative \rightarrow penalty.

if positive \rightarrow no penalty.

0-1 loss: $\ell_{0-1}(z) = \begin{cases} 1 & \text{if } z \leq 0 \quad (\text{mistake} \rightarrow \text{penalty}) \\ 0 & \text{if } z > 0 \quad (\text{no mistake} \rightarrow \text{no penalty}) \end{cases}$

for classification

The optimization problem associated with classification:

$$\min_w \frac{1}{m} \sum_{i=1}^m \ell_{0-1}(y^i \cdot w^T x^i)$$

Want to solve this.

When does perceptron find w with small loss?

Recall the perceptron required $\exists w$ s.t. $\forall x$:

$$\underset{\substack{\uparrow \\ y \in \{-1, +1\}}}{y} \cdot w^T x > \rho \quad \Rightarrow \quad \begin{array}{l} \text{convergence} \\ \# \text{ of mistake} \end{array} < \frac{1}{\rho^2}$$

What if there is no margin?

Also, there might not exist a w s.t. $\text{sign}(w^T x^i) = y^i \forall i$,
e.g. they are not linearly separable, e.g. there is not
a half space that correctly classifies all the points.

Intuitively, we can still try to $\min_w \left[\frac{1}{m} \sum_{i=1}^m \mathcal{L}_{0-1}(y^i w^T x^i) \right]$

\Rightarrow To find a half-space that maximizes the # of correct
labels, even though it cannot fully separate correct from
incorrect,

what is the computational complexity of this optimization
problem?

Unfortunately, \mathcal{L}_{0-1} is neither convex nor differentiable.

This problem is NP hard, e.g. unlikely to have
polynomial time solution.

Also it is "agnostically learning a half-space"

Summary :

Regression \longrightarrow convex loss function.

Classification \longrightarrow non convex loss (0-1). (Bad news),

Idea: Relax the 0-1 loss to a different nicer loss:
"surrogate loss", related to 0-1 loss but convex.
would get solutions that are close or high-quality.

Let's introduce a few losses: logistic, hinge, exponential.

$$\bullet \text{Logistic}(z) = \log(1 + e^{-z})$$

$$\ell_{\text{logistic}}(y^i \cdot w^T x^i) = \log(1 + e^{-(y^i \cdot w^T x^i)})$$

if $\underbrace{y^i w^T x^i}_{\text{margin}} \ll 0 \Rightarrow \ell_{\text{logistic}}(y^i \cdot w^T x^i)$ is large.

" $\gg 0 \Rightarrow \ell_{\text{logistic}}(y^i \cdot w^T x^i)$ is small,
(moves to 0).

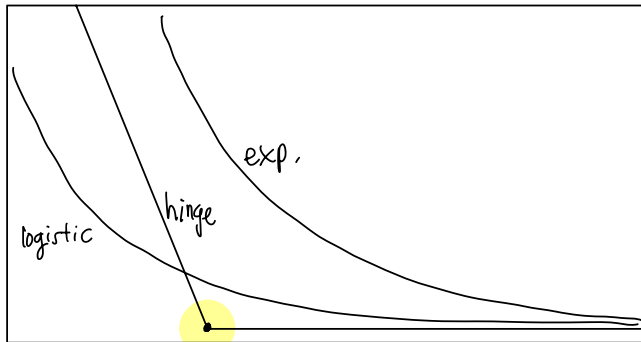
$$\bullet \text{hinge}(z) = \max\{1 - z, 0\}$$

$$\ell_{\text{hinge}}(y^i \cdot w^T x^i)$$

\Rightarrow large when $y^i \cdot w^T x^i$ is negative.

small when $y^i \cdot w^T x^i$ is positive.

$$\bullet \text{exp} = e^{-z},$$



$$(z = yw^T x) \rightarrow "x=1"$$

Logistic Loss Optimization :

average logistic loss

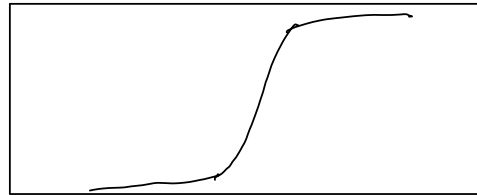
$$\downarrow$$

$$L(w) \equiv \frac{1}{m} \sum_{i=1}^m \underbrace{\log(1 + \exp(-y^i \cdot w^T x^i))}_{\text{logistic}}$$

Goal : $\min_w [L(w)]$

Enter the sigmoid function :

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



$$\text{As } z \rightarrow \text{large} \Rightarrow g(z) \rightarrow 1$$

$$\text{As } z \rightarrow \text{small} \Rightarrow g(z) \rightarrow 0$$

fact : $g(z) + g(-z) = 1$. Proof :

$$\frac{1}{1 + e^{-z}} + \frac{1}{1 + e^z} \Rightarrow \frac{e^z}{e^z} \left(\frac{1}{1 + e^z} \right) + \frac{1}{1 + e^z} \Rightarrow \frac{e^z}{e^z + 1} + \frac{1}{1 + e^z} = 1$$

$$E[Y|X] = g(y \cdot w^T X) \quad \text{for some } w. \quad Y \in \{-1, +1\}$$

$$\Rightarrow \Pr[Y=1|X] = g(w^T X) = g(w^T x)$$

Given X :

- if $w^T X$ is large: $\Rightarrow \Pr[Y=1|X] = \text{large}$
 - if $w^T X$ is small or neg: $\Rightarrow \Pr[Y=1|X] = \text{small}$.
- soft relaxation of half-space.

$$\Pr[Y=y^i | x^i; w] = g(y^i \cdot w^T x^i), \text{ where } g(x) = \frac{1}{1+e^{-x}}$$

"model for logistic regression"

Given a training set $S = \{(x^1, y^1), \dots, (x^m, y^m)\}$, what is the most likely w given the training set? Use **MLE**.

$$\text{Likelihood}(w) = \prod_{i=1}^m \Pr(Y=y^i | x^i, w) = \prod_{i=1}^m g(y^i \cdot w^T x^i)$$

max w \nearrow

$$\begin{aligned} \log\text{-likelihood} &= \sum_{i=1}^m \log g(y^i \cdot w^T x^i) = \sum_{i=1}^m \log \left(\frac{1}{1 + \exp(-y^i w^T x^i)} \right) \\ &= - \sum_{i=1}^m \underbrace{\log(1 + \exp(-y^i \cdot w^T x^i))}_{\text{logistic loss}} = \underbrace{-m \cdot L(w)}_{\text{average logistic loss} = \frac{1}{m} \sum l_{\text{logistic}}} \end{aligned}$$

Log-likelihood is maximized when $L(w)$ is minimized, due to $-m$.

Now our goal is to minimize the logistic loss $L(w)$.

Idea: Run gradient descent on logistic loss.

This is the algorithm for performing logistic regression.

Note: There is no closed form, hence need for GD.

Let's compute the gradient of $l(w)$:

$$l_{\text{logistic}}(z) \equiv \log(1 + e^{-z})$$

$$1) \quad l'_{\text{logistic}}(z) = \frac{1}{1 + e^{-z}} \cdot -e^{-z} = \frac{-1}{1 + e^{+z}} = -g(-z)$$

$$2) \text{ compute } \frac{\partial l_{\text{logistic}}(y \cdot w^T x)}{\partial w_k} = \underbrace{-g(-y \cdot w^T x)}_{\substack{\uparrow \text{sigmoid function.}}} \cdot \underbrace{y \cdot x^k}_{\text{chain rule } \frac{\partial z}{\partial w_k}}$$

with this formula we can directly apply gradient descent;
this precisely tells us how to find **max-likelihood w** .

What happens if we have multiple labels for y ?

What if $y \in \{1, \dots, k\}$?

use **multinomial logistic regression**: w^1, \dots, w^k weight vectors

$$P[y=i|x] = \frac{e^{w^i \cdot x}}{\sum_{i=1}^k e^{(w^i)^T x}} \quad \text{or } k-1$$

$$P[y=1|x] \propto e^{(w^1)^T x}$$

$$P[y=i|x] \propto e^{(w^i)^T x}$$

$$P[y=k] = 1 - \sum_{i=1}^{k-1} P[y=i]$$

What is the associated loss? cross-entropy loss,
Generalization of logistic loss.

(Imagine y is a vector of length K with a 1 in the j^{th} position if correct label is j).

one-hot encoding of labels

Let's say our guess for the probability y has label i is p_i : $-\sum_{i=1}^k y_i \log(p_i)$ (this is cross-entropy loss).

softmax \rightarrow turns real-values into probabilities.

Example:

$w^T x$ via $\text{sigmoid}(w^T x) \xrightarrow{\text{map}} [0, 1]$ probability space.

$\underbrace{(z_1, \dots, z_k)}_{k \text{ coordinates}} \rightarrow \underbrace{\left(\frac{e^{z_1}}{Z}, \frac{e^{z_2}}{Z}, \dots, \frac{e^{z_k}}{Z} \right)}_{\text{probability space; sums to 1,}}$

$$\Rightarrow Z = \sum_{i=1}^k e^{z_i}$$