

KEY CONCEPTS:

Dec bound
Loss fn
S.G.D.

feature vector.

CATEGORICAL y : $\vec{x} \in \mathbb{R}^d \rightarrow \{1, \dots, K\}$ or $\{0, 1\}$, $\{-1, 1\}$. (d = dimension X_1, \dots, X_d).

Remember: $\hat{y} = f(\vec{x}, \vec{w}) = \vec{w}^T \vec{x}$ (we still have a BIAS).

Binary decision:

$$\hat{y} = f(\vec{x}, \vec{w}) = \text{sign}(\vec{w}^T \vec{x}) = \begin{cases} 1 & \text{if } \vec{w}^T \vec{x} \geq 0 \\ -1 & \text{if } \vec{w}^T \vec{x} < 0 \end{cases} \text{ threshold.}$$

$X_1 \ X_2 \ \dots$
1 1 ...

$$\vec{X} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \begin{matrix} \text{-data pt} \\ \text{-data pt} \\ \cdot \\ \cdot \\ \cdot \end{matrix}$$

$\vec{w}^T \vec{x} = 0$ is HYPERPLANE.
ORTHOG TO \vec{w} .

$\vec{w}^T \vec{x} + w_0 = 0$ shifts by w_0 (bias).

LINEARLY SEPARABLE: separate classes by hyperplane,
($d-1$ plane).

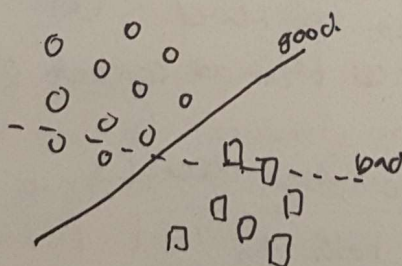
some issues that make it not linsep:

- model too simple (polynomial as polynomial)
- bad featur.
- noise

Should we have perfect separation in training data?
Just b/c we have perfect sept. in training doesn't mean we have it in test.

Need to find \vec{w} (direction), w_0 (location) of boundary.

What is a "good" boundary?



→ We need a Loss FN.

Two options:

- NATURAL ZERO-ONE LOSS

$$l_{0-1}(\hat{y}, y) = \begin{cases} 1 & \text{if } y \neq \hat{y} \\ 0 & \text{if } y = \hat{y} \end{cases}$$

BUT HARD TO MINIMIZE--

nonconvex, noncontinuous piecewise constant.

NP-H
to
optimize

- ASYMMETRIC BINARY LOSS:

"two types of mistakes"

$$l_{ABL}(\hat{y}, y) = \begin{cases} \alpha & \text{if } y=0 \wedge \hat{y}=1 \\ \beta & \text{if } y=1 \wedge \hat{y}=0 \\ 0 & \text{if } y=\hat{y} \end{cases}$$

} if repercussions
for 2 types of
mistakes are
different.

- USE A SURROGATE
LOSS INSTEAD, $\hat{\ell}$.

Good surrog. loss

- Easy to optimize

- Representative:

low s. loss means low orig. loss

$l_2(y, \hat{y}) = (y - \hat{y})^2$ is easy to optimize, but
not on representation.
(very harsh on inliers).

WE USE METRICS to evaluate loss functions

- METRICS:

- ACCURACY: % of correct predictions

noninformative when considering RARE OUTCOMES.

TP = true pos

FN = False neg

FP = False pos

Recall: relevant instances retrieved (R)

Precision: retrieval instances that are correct (P)

F₁ SCORE: HARMONIC MEAN OF PRECISION AND RECALL

$$F_1 = 2 \frac{PR}{P+R} \quad \left. \vphantom{F_1} \right\} \begin{array}{l} \text{need both to} \\ \text{be "good" to} \\ \text{work.} \end{array}$$

PRECISION-RECALL CURVE.

tradeoff btwn precision, recall via decision threshold.

OPTIMIZATION.

How do we find w once we have a loss fn l ?

WE USE GRADIENT DESCENT not b/c it is THE BEST, but b/c it is EASY TO SCALE. It is best for SOLVING BG PROBLEMS.

$$w = \frac{1}{N} \sum_{i=1}^N l(x_i, f(x_i, w)).$$

GRADIENT DESCENT.

- Initialize w_0 .
- Compute gradient (tangent dir).
- Take a step.
- Keep updating

$$w_{t+1} = w_t - \underbrace{\lambda_t}_{\text{LEARNING RATE}} \nabla_w L(w_t).$$

λ :

LARGE $\lambda \rightarrow$ overshoot

SMALL $\lambda \rightarrow$ long time.

line search (searching min along gradient line) too slow.

standard: decay λ as learning progresses

THERE IS A LOT OF THEORY ON CONVEX OPTIMIZATION. (ask Refaya for the literature).

Comp. cost of computing $\nabla_w L(w_t)$?

$$L(w) = \frac{1}{N} \sum_{i=1}^N l(x_i, f(x_i, w)) = \text{grows LINEARLY in } N.$$

HUGE DATASET \rightarrow large cost for small update

SOLN TO ABOVE:

Pick one pt randomly instead of gradient over all.

Find gradient at that

THIS IS STOCHASTIC GRADIENT DESCENT

Pick datum:

$$g_t = \nabla l(x_j, f(x_j, w)).$$

Compute

Theoretically works, but practically VERY NOISY.

INSTEAD OF ONE PT (S.G.D) vs ALL PTS (G.D).

pick a "small subset" MINIBATCH.

tradeoff: high batch accuracy \propto small batch runtime.

Instead of mapping X to $\{0, 1\}$,
 we now map to $[0, 1]$ and we find $P(y=1|x)$.
 (STILL CLASSIFICATION)

Need to SQUASH wTx into $[0, 1]$ & $P(y=1|x)$.
 What about $P(y=-1|x)$? (Prob of other class).
 $P(y=-1|x) = 1 - P(y=1|x)$.

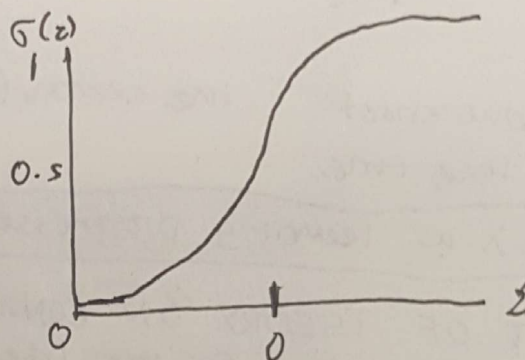
How to CHOOSE THE LABEL?
 Pick the MOST PROBABLE.

BENEFITS:

- MODELS UNCERTAINTY (in LIMITED way)
- Can use Pr for decision making
- Can use Pr for prob/stat opt. methods.

SQUASHING FXN:

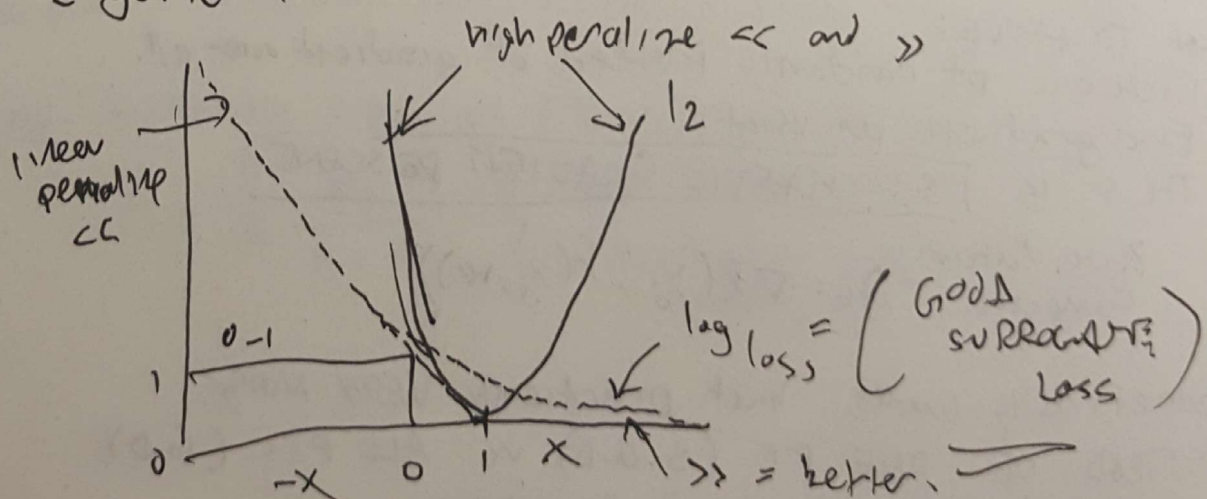
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Modifying w and w_0
 changes the shape!

w = "flatness/width" \rightarrow SLOPE / sharpness of $0 \rightarrow 1$ convergence
 w_0 = "location".

High $w \approx$ step fn.
 Low $w \approx$ gentle increase.



logistic regression: optimize surrogate loss.