

$$= \sum_{i=1}^N \log \pi_i + \log \left[\text{Gaussian dist} \right]$$

$$+ \sum_{i=1}^N \log \pi_2 + \log \left[\quad \right]$$

$$\therefore N_1 \log \pi_1 + \sum_{y_i=1}^N \left[\quad \right] + \sum_{y_i=2}^N \left[\quad \right]$$

$$+ N_2 \log \pi_2 - \frac{N \log (\Sigma)}{2} + \text{constant}$$

- Maximize LL wrt using $\pi_1, \pi_2, \mu_1, \mu_2, \Sigma$ subject to $\pi_1 + \pi_2 = 1$

Find gradient wrt each parameter. Make it zero.

This is valid only if LL is concave

Contd on next pag

10/8

* Convex Functions and Optimization

$$\vec{\theta} = \theta_1, \theta_2, \dots, \theta_r \quad \text{Find } \min F(\vec{\theta})$$

$$\theta_i \in \mathbb{R}, \vec{\theta} \in \mathbb{R}^r$$

$F(a) \in \{\text{convex, concave, neither}\}$

1. $F(a)$ is convex.

* If $F(a)$ is convex, then $-F(a)$ is concave.

Definition $F(\theta)$ is convex iff $\theta_1^*, \theta_2^* \in \mathbb{R}^r$, $\forall \lambda \in [0, 1]$,

$$\lambda F(\theta_1) + (1-\lambda) F(\theta_2) \geq F(\lambda \theta_1 + (1-\lambda) \theta_2)$$

Chord lies above curve

- For local minimas of differentiable $F(\theta)$.

$$\nabla F(\theta^*) = 0$$

Definition Gradient :- $\nabla F(\vec{\theta}) = \begin{bmatrix} \frac{\partial F}{\partial \theta_1} \\ \vdots \\ \frac{\partial F}{\partial \theta_r} \end{bmatrix}$

Continued from prev page.

$$F(\theta) = N_1 \log \pi_1 + \sum_{i=2}^N \frac{1}{2} \frac{(x_{1i} - \mu_{1i})^2}{\sigma^2} + \sum_{i=2}^N \frac{1}{2} \frac{(x_{2i} - \mu_{2i})^2}{\sigma^2} + N_2 \log \pi_2 - N \log \sigma + \text{const}$$

Maximize w.r.t $\pi_1, \pi_2, \mu_{11}, \mu_{21}, \sigma$ such that $\pi_1 + \pi_2 = 1$

$$\text{Write } \pi_2 = 1 - \pi_1$$

$$\therefore \nabla LL(\theta) = \begin{bmatrix} \frac{N_1}{\pi_1} - \frac{N_2}{1-\pi_1} \\ \sum \frac{(x_{1i} - \mu_{1i})}{\sigma^2} \\ \sum \frac{(x_{2i} - \mu_{2i})}{\sigma^2} \\ \sum_{i=3}^N \frac{(x_{1i} - \mu_{1i})^2}{\sigma^3} + \sum_{i=3}^N \frac{(x_{2i} - \mu_{2i})^2}{\sigma^3} - N \end{bmatrix} = \vec{0}$$

$$\pi_1 = \frac{N_1}{N_1 + N_2}, \quad \mu_{11} = \frac{\sum_{i=1}^N x_{1i}}{N_1}, \quad \mu_{21} = \frac{\sum_{i=2}^N x_{2i}}{N_2}$$

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_{1i} - \mu_{11})^2 + \sum_{i=2}^N (x_{2i} - \mu_{21})^2}{N}}$$

III

CONDITIONAL PROBABILISTIC CLASSIFICATION

Logistic Regression.

Assume $P(y|x) \sim \text{Bernoulli}(\alpha)$, where $y \in \{1, 2\}$

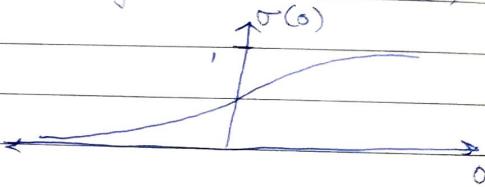
- Parameter of Bernoulli depends on inputs.

General Trick Compute a function g of x_1, \dots, x_d , such that

$$g(x) \in [-\infty, \infty]$$

Squash this $g(x)$ to be between $(0, 1)$ so as to look like a Bernoulli parameter using the sigmoid function $\sigma(x)$

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



- Define $P(y=1|x) = \sigma(g(x)) = \frac{1}{1 + e^{-g(x)}}$

$$P(y=2|x) = 1 - P(y=1|x) = \frac{1}{1 + e^{g(x)}}$$

Define $y=2 \rightarrow y=-1$.

$$P(y|x) = \frac{1}{1 + e^{-yg(x)}} \quad y \in \{-1, +1\}$$

$g(x)$ is learned from data.

→ Linear Logistic Regression.

$$g(x) = w^T x + w_0 \text{ where.}$$

$$w \in \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

for $w_i \in \mathbb{R}$

$$g(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 \dots w_0$$

$$P(y|x) = \frac{1}{1 + e^{-y(w_0x_0 + \dots + w_dx_d)}}$$

for $y = \{-1, 1\}$

* Is $\sum P(y_i|x) = 1$?

$w_0 \text{ to } w_d$: The $d+1$ parameters of the classifier that will be learned from training data.

- $D \equiv \{(x_1, y_1), \dots, (x_n, y_n)\}$

Find w_p so as to maximize the conditional likelihood of D .

- Au contraire, for generative classifiers, we fit joint likelihood

Conditional

$$\max_{\theta} \sum_{i=1}^N \log P_r(y_i | x_i, \theta \equiv w)$$

Generative

$$\max_{\theta} \sum_{i=1}^N \log P_r(x_i, y_i | \theta)$$

* Conditional is better than generative because it does not make unreasonable assumptions (especially for high-dimensional data)

- Maximize $\sum_{i=1}^N \log P(y_i | x_i, \theta)$ wrt $\theta = \{w_0, \dots, w_d\}$

$$= \max_w \sum_{i=1}^N \log \frac{1}{1 + e^{-y_i(w^T x_i + w_0)}}$$

* Sigmoid function is concave.

$$\sigma''(w) = \frac{-e^w}{(1+e^w)^2} < 0$$

* If $F(\alpha)$ is concave in α and $\alpha = w^T x + w_0$, then F is also concave in (w_0, \dots, w_d)

$\max LL(w, w_0, D)$ can be solved by $\nabla LL = 0$

$$LL = \sum_{i=1}^N -\log \left(1 + e^{-y_i(w^T x_i + w_0)} \right)$$

$$\nabla LL = - \left[\sum_{i=1}^N \frac{e^{-y_i(w^T x_i + w_0)}}{1 + e^{-y_i(w^T x_i + w_0)}} - y_i \right] \cdot (-y_i x_i)$$

e.g. $d=1$, $N=2$, $D = [(1, 1), (-1, -1)]$

$$\nabla LL = - \left[\begin{array}{c} \frac{e^{-(w_1 + w_0)}}{1 + e^{-(w_1 + w_0)}} (-1) + \frac{e^{(-w_1 + w_0)}}{1 + e^{(-w_1 + w_0)}} (1) \\ \frac{e^{-(w_1 + w_0)}}{1 + e^{-(w_1 + w_0)}} (-1) + \frac{e^{(-w_1 + w_0)}}{1 + e^{(-w_1 + w_0)}} (-1) \end{array} \right] = 0$$

Find w_1 & w_0

17/2 is difficult to solve

- Q. Use 'Gradient Ascent Algorithm', which is an iterative algorithm
 - start with a guess. Try to converge to the solution.

- Find the decision boundary by solving the following
- $$\log P(y=1 | x) > \log P(y=-1 | x)$$

$$\log \left(\frac{1}{1 + e^{-(w^T x + w_0)}} \right) > \log \left(1 - \frac{1}{1 + e^{-(w^T x + w_0)}} \right)$$

$$= \log \left(\frac{1}{1 + e^{-(w^T x + w_0)}} \right)$$

$$\therefore 0 > -(\mathbf{w}^T \mathbf{x} + w_0)$$

$$\mathbf{w}^T \mathbf{x} + w_0 > 0$$

- Logistic Regression also gives a linear decision boundary.
- Logistic Regression gives better accuracy than LDA.

2A/8

- Probabilistic classifiers are easy to extend to more than two classes

→ For more than two classes,

$$y \in \{1, 2, 3, \dots, K\}$$

- ∵ y follows a Multinomial distribution.

$$y \sim \text{Multinomial}(p_1, p_2, \dots, p_K)$$

$$\text{such that } 0 < p_i < 1, \quad \sum p_i = 1$$

- $P(y|x) \sim \text{Multinomial}(p_{1x}, p_{2x}, \dots, p_{Kx})$

$$\text{such that } 0 < p_{ix} < 1, \quad \sum p_{ix} = 1$$

- Parameters are dependent on value of x .

$$P(y=j|x) = P_{jx} = \frac{e^{w_j x + w_{j0}}}{\sum e^{w_y x + w_{y0}}}$$

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

Use likelihood estimator.

$$\text{No. of parameters} = k(d+1) \quad \text{No. of } W_{ij}$$

- * One of the optimal solutions is that W_s of $+y$ are negative of W_s of $-y$ (like in the case of binary classification, where we directly used only a single set of W_s)

- Likelihood estimator = $LL(D, w)$
 $= \sum \log [P(y_i | w, \vec{x}_i)]$
 $= \sum_{i=1}^N \left(w_{y_i}^\top \vec{x}_i + w_{y_i^0} - \log \left(\sum e^{w_{y_i}^\top \vec{x}_i + w_{y_i^0}} \right) \right)$

This function is concave in w .

- * If $F(z)$ is concave in z , then $F(w^\top x)$ is concave in w ,
 ↪ with x as a constant and w is multidimensional

- Maximize $LL(D, w)$ wrt w .
 - Not trivial
 - Use iterative algorithm

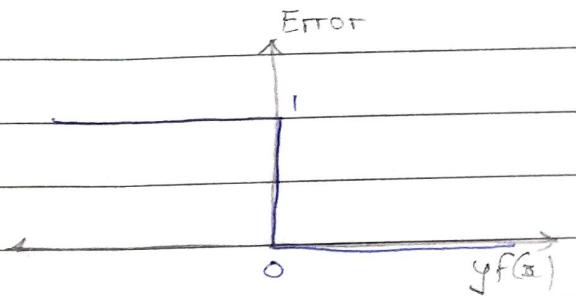
IV Loss REGULARIZATION FAMILY OF CLASSIFIERS

eg. $k = 2$, $y_i \in \{-1, 1\}$

$$f(x) = w \cdot x + w_0$$

\uparrow
bias

- Classifier := sign of $f(x)$
- Error = 1 if $y f(x) < 0$
= 0 otherwise



Training $D = \{(x_i, y_i) \dots\}$ $1 \leq i \leq N$

TFT = Ideal value of w

$$= \operatorname{argmin} \left(\sum \text{error}(y_i, f(x_i | w)) \right)$$

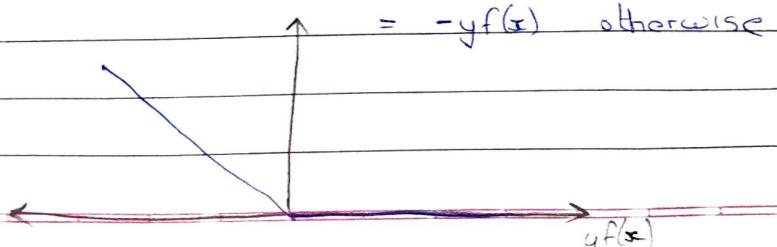
$$= \operatorname{argmin} \left(\sum \text{error}(y_p, (w \cdot x_p + w_0)) \right)$$

- Error ($y f(x)$) is neither concave nor convex.
 - Line joining any two points does not lie entirely below or above the function.
 - Cannot use concave optimization problem
 - Define new loss functions

→ Perception Loss Function.

$$\text{ploss}(y f(x)) = 0 \text{ if } y f(x) \geq 0$$

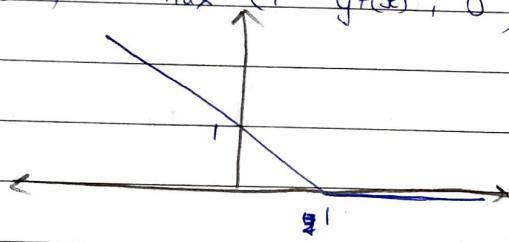
$$= -y f(x) \text{ otherwise}$$



- ∴ ploss is convex and continuous
- ∴ ploss is not differentiable
- ∴ For this classifier, margin of separation is very thin
 - $0.00001 \rightarrow y = 1$
 - $-0.00001 \rightarrow y = -1$
 - We want decision boundaries that pass through regions which \nexists have low density of points

→ Hinge Loss

- $\text{Hinge}(yf(x)) = \begin{cases} 1 - yf(x) & \text{if } yf(x) \leq 1 \\ 0 & \text{otherwise.} \end{cases}$
- $\text{Hinge}(yf(x)) = \max(1 - yf(x), 0)$



- ∴ The shift of 1 unit boosts generalizability.
Error is zero only if $yf(x)$ is significantly more than zero

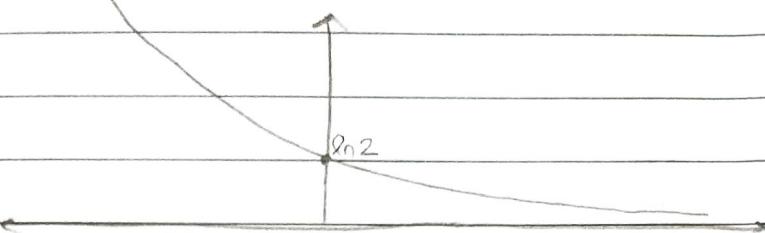
→ Square Hinge Loss

$$= \text{Hinge}^2(yf(x))$$

PAGE NO.	
DATE	/ /

→ Logistic Loss

$$\text{logloss}(y f(x)) = \log(1 + e^{-y f(x)})$$



$\approx -z$ for $z \rightarrow -\infty$

∴ - Convex in $y f(x)$

- Differentiable

- Gives a good margin. - does not overfit.

9/8

A]

Regularizers

- Prevent overfitting

e.g. 1-dimensional data :- $D \equiv \{(x_1, y_1), \dots, (x_n, y_n)\}$ $d=1$

Use a linear function. If it doesn't overfit.

- If above data is synthetically expanded to g -dimensional data,

$$D' \equiv \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$$

where, say, $\vec{x}_i = (x_i, x_i^2, x_i^3, \dots, x_i^g)$

~~Linear~~ Even linear functions will overfit on this g -dimensional space.

ditch
18

- Example of Regularizer :- $R(\omega) = \left(\sum_{j=1}^d |\omega_j|^\gamma \right)$

Training objective :- Minimize - $C \left(\sum_1^N \text{loss}(y_i f(x_i)) \right) + R(\omega)$
wrt ω, ω_0

- Note :-

- Bias term (ω_0) is not regularized.
- C can be chosen depending on how much weight you want to give to $R(\omega)$.

- Value of γ :-

• $\gamma = 0$:- Only ^{no. of} non-zero w_j 's are minimized. L0

• $\gamma = 1$ L1

• $\gamma = 2$ L2

- Used most often.

$$R(\omega) = \sum_1^d \omega_j^2$$

- Like the loss function, regularizer is also a convex function in w_j . ☺

B]

Algorithms for optimizing convex functions,

Assumption :- Solving $\nabla F(\omega) = 0$ is not easy.

Objective function :- $F(\omega) = C \left(\sum \text{loss}(y f(x)) \right) + R(\omega)$

- Four types of algorithms :- 0th, 1st, 2nd, semi 2nd

- 1 Zeroth order :- Gradient Computation
 Applicable to non-differentiable function.
- 2 First order :- Gradient Computation
 Applicable to functions that are not twice differentiable (eg - Square Hinge loss.)
- 3 Second order :- Hessian Computation
 ↙ analogous to second differentiation in multidimensional data.
- 4 Semi-second order :- Approximate the Hessian.

1/8

A Zeroth Order :- Only $f(w)$ is known

→ Calculate $f(w)$ for many many points, $\Delta w = \epsilon$ apart, and return the minimum $f(w)$ found.

Tolerance of ϵ is applicable

∴ Can be applied for any function, even non-convex functions

∴ Duh.

→ Dichotomous Binary Search

a and b are on opposite sides of minima.

while ($|b-a| \geq \epsilon$):

Pick two points w_1 and w_2 , $\frac{\epsilon}{2}$ apart, in the middle of a and b

IF ($F(w_1) > F(w_2)$): $a = w_1$

Else $b = w_2$

$$\text{No. of iterations} = T_2$$

function calls

$$\text{If } T = \text{no. of evaluations}, (b-a)(0.5)^{\frac{T}{2}} = \epsilon$$

$$\Rightarrow T = \frac{2 \ln \epsilon}{\ln(0.5)}$$

- This algorithm surely terminates, and gives correct minimum
- Cannot be used in higher dimensions

B

First order methods

$f(w)$ and $\text{grad } f(w)$ are known.

→ 1D function data :- Modified Bisection Algorithm

while $|b-a| \geq \epsilon$..

Choose $w \in (a, b)$

If $f'(w) = 0$:- return w

> 0 :- $b = w$

< 0 :- $a = w$

→ First order methods for higher dimensional data

- Taylor Series expansion

$$F(w) = F(w^0) + \nabla F(w^0)^T (w - w^0) + |w - w^0| \varphi(w^0, w)$$

such that $\varphi(w^0, w) \rightarrow 0$ as $w^0 \rightarrow w$

Hence, every differentiable function can be linearly approximated in the neighbourhood of any point

$$\tilde{F}(w) = F(w^0) + \nabla F(w^0)^T (w - w^0)$$

Write $w = \lambda \vec{e} + w^0$, where \vec{e} is a unit vector in d dimensions and $\lambda \in \mathbb{R}$, $\lambda \geq 0$, λ is small.

$$\tilde{F}(w^0 + \lambda \vec{e}) = F(w^0) + \lambda \underbrace{\nabla F(w^0)^T \cdot \vec{e}}_{\leq 0}$$

Move in direction opposite to gradient:- Choose \vec{e} such that $\nabla F(w^0)^T \cdot \vec{e} \leq 0$

Then you will get $F(w) \leq F(w^0)$

Algorithm $w^0 = \text{arbitrary initial value}$

while $(|\nabla F(w^t)| \geq \epsilon)$ {

i) Find \vec{e} s.t. $\nabla F(w^t) \cdot \vec{e} \leq 0$

ii) $\lambda^* = \arg \min [F(w^t + \lambda \vec{e})]$... ID convex optimization

iii) $w^{t+1} = \lambda \vec{e} + w^t$.

→ Gradient Descent Algorithm:-

- Choose $\vec{e} = -\nabla F(w^0)$

- Pick λ as any small positive value. (avoids calculation of λ^*)

- Applying for logistic loss minimization,

$$F(w) = \sum_{i=1}^N \log \left(1 + e^{-w x_i y_i} \right)$$

$$\nabla F(w) = \sum_{i=1}^N \frac{e^{-w x_i y_i}}{1 + e^{-w x_i y_i}} (-x_i y_i) \text{ --- vector}$$

→ Off-the-shelf optimization algorithm :- 'LBFGS'

- The loop is implemented by the algorithm.

You only need to evaluate $F(w)$ and $\nabla F(w)$ at any given w .

→ Stochastic Gradient Descent Algorithm.

Instead of calculating precise gradient, calculate gradient using only a subset of data points

This subset is called a 'batch'

- Makes execution time lower.

C

Second Order Methods

- For a twice differential function, Hessian is a $d \times d$ matrix where d is the dimensionality of data

$$H_{ij} = \frac{\partial^2 F}{\partial w_i \partial w_j}$$

- For any twice differentiable function, second order Taylor expansion:

$$F(w) = F(w^0) + \nabla(F(w^0))^T (w - w^0) + \frac{1}{2} (w - w^0)^T H_F(w^0) (w - w^0) + \|w - w^0\|^2 [\alpha(w^0, w - w^0)]$$

where $\alpha \rightarrow 0$ as $w^0 \rightarrow w$