

Naiive Bayes and Logistic Regression

Nimisha Roy
Georgia Tech

Outline

Naive Bayes
Logistic Regression

- Generative and Discriminative Classification
- The Logistic Regression Model
- Understanding the Objective Function
- Gradient Descent for Parameter Learning
- Multiclass Logistic Regression



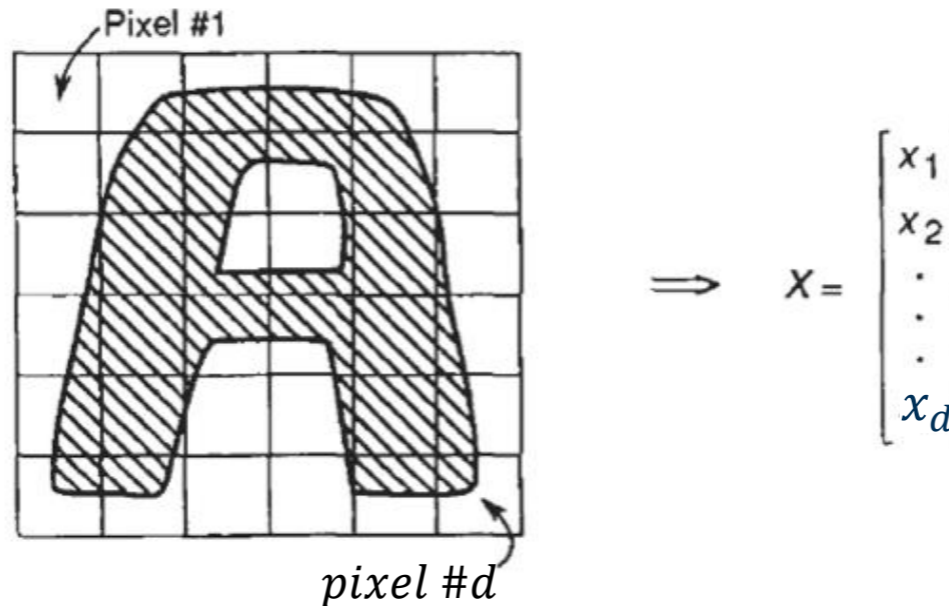
$$P(y=1|x)$$

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} = \frac{P(x,y)}{P(x)}$$

$P(y|x)$ is computed directly

Classification

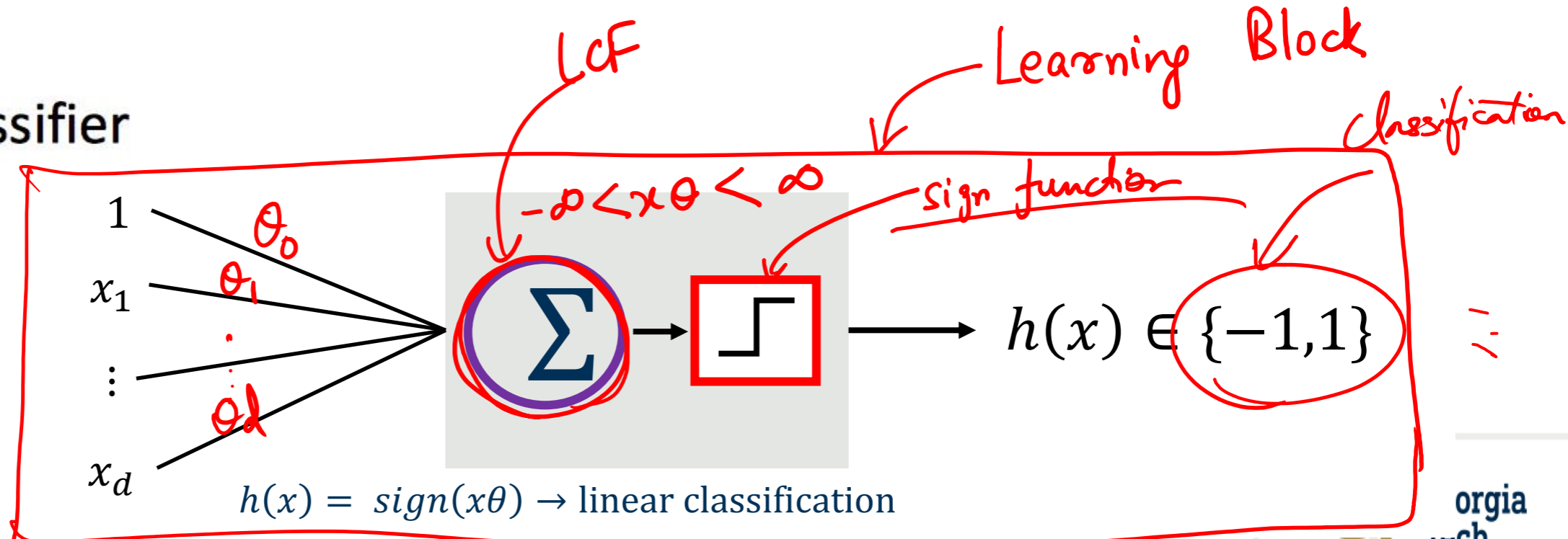
- Represent the data



- A label is provided for each data point, eg., $y \in \{-1, +1\}$

A (with arrow pointing to +1)
not A (with arrow pointing to -1)

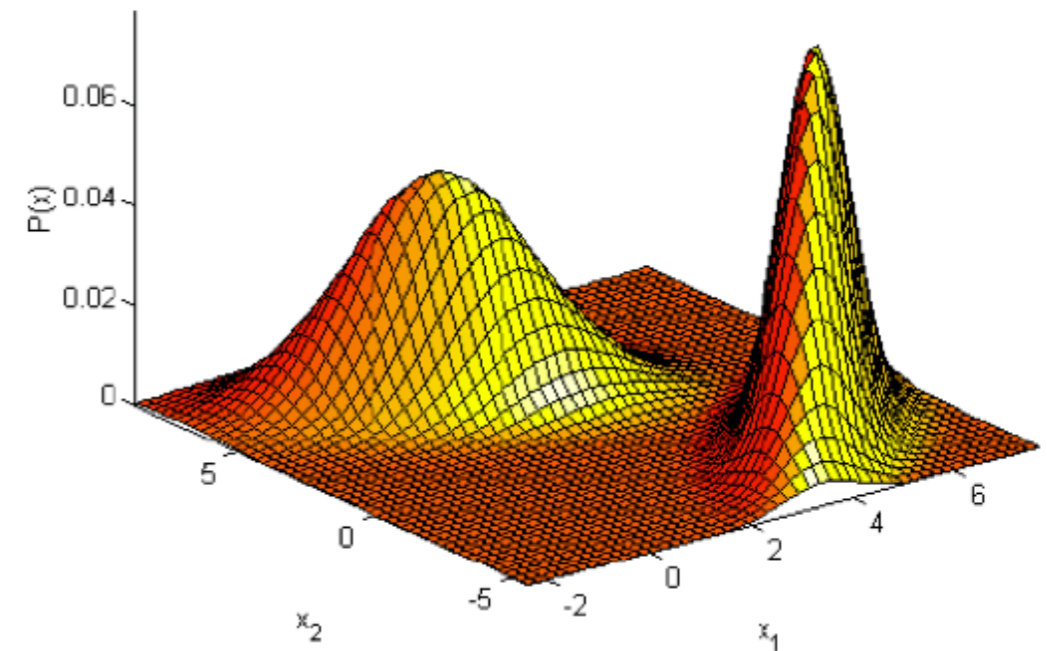
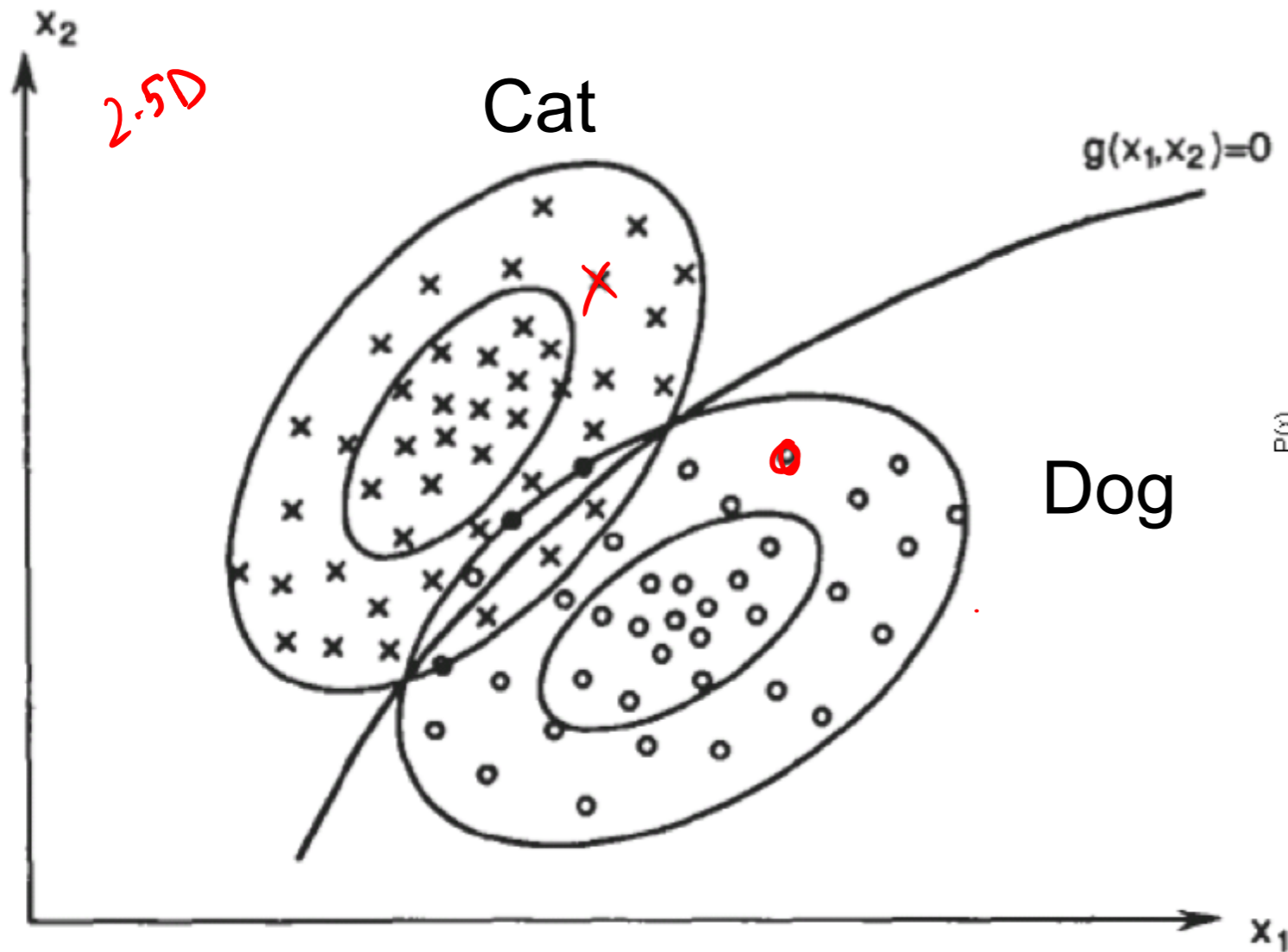
- Classifier



~~$\{i\}$~~

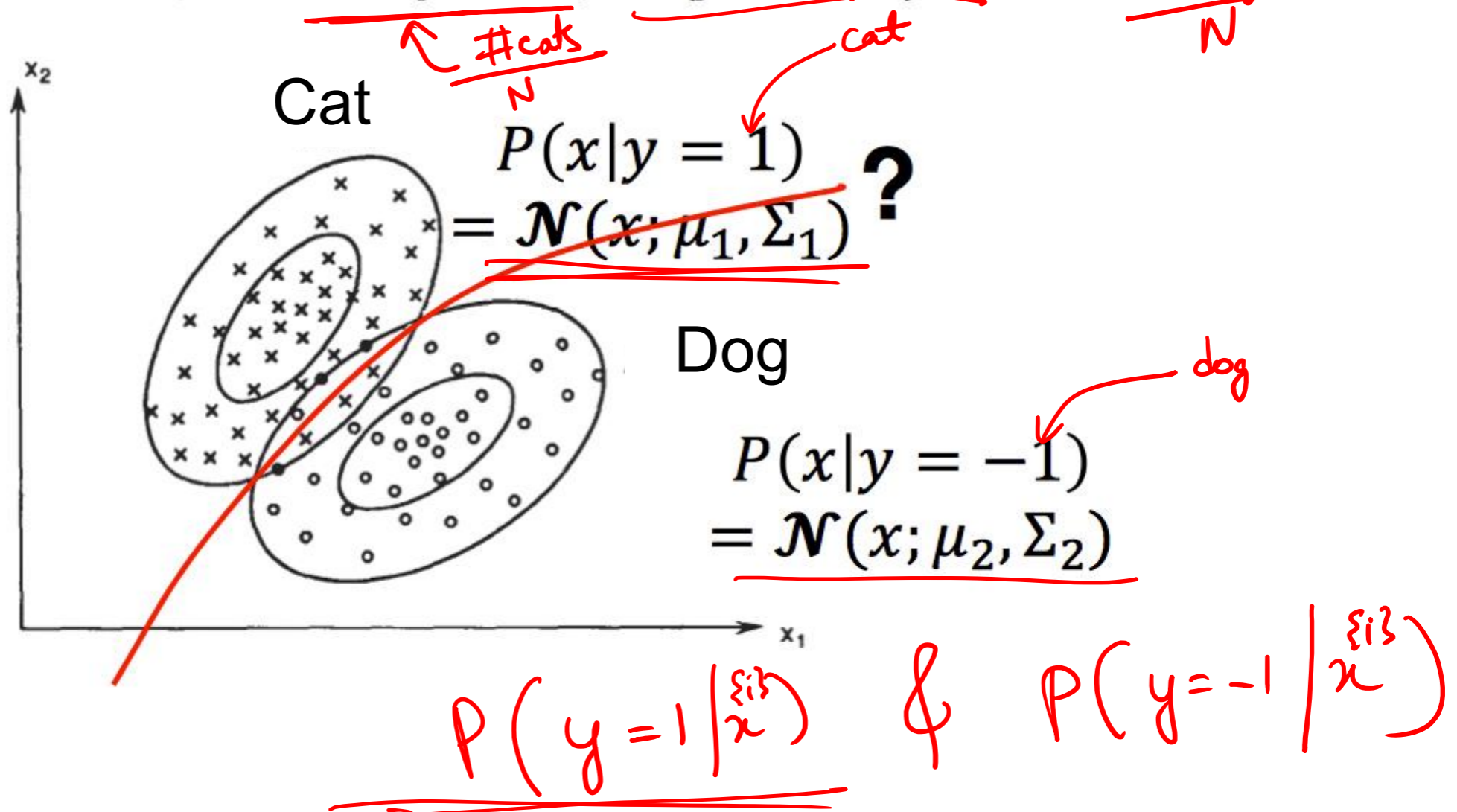
Decision Making: Dividing the Feature Space

- Distributions of sample from normal (positive class) and abnormal (negative class) tissues



How to Determine the Decision Boundary?

- Given class conditional distribution: $P(x|y = 1), P(x|y = -1)$, and class prior: $P(y = 1), P(y = -1)$



Bayes Decision Rule

$P(y=1|x)$ vs. $P(y=-1|x)$

$$\underbrace{P(y|x)}_{\text{posterior}} = \frac{\overbrace{P(x|y)P(y)}^{\text{likelihood} \times \text{Prior}}}{\underbrace{P(x)}_{\text{normalization constant}}} = \frac{P(x, y)}{\sum_z P(x, y)}$$

Prior: $P(y)$

Likelihood (class conditional distribution : $p(x|y) = \mathcal{N}(x|\mu_y, \Sigma_y)$

$$\text{Posterior: } P(y|x) = \frac{P(y)\mathcal{N}(x|\mu_y, \Sigma_y)}{\sum_y P(y)\mathcal{N}(x|\mu_y, \Sigma_y)}$$

$$P(y = \pm 1 | x^{\xi i}) = \frac{P(x^{\xi i} | y = \pm 1) \cdot P(y = \pm 1)}{P(x^{\xi i})}$$

$$= \frac{P(x_1^{\xi i}, x_2^{\xi i}, x_3^{\xi i}, \dots, x_d^{\xi i} | y = \pm 1) \cdot P(y = \pm 1)}{P(x^{\xi i})}$$

$$= \frac{1}{2\pi |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right) \cdot P(y = \pm 1)$$

Naive Bayes \rightarrow Assume Conditional independence between features

$$= \frac{P(x_1^{\xi i} | y = \pm 1) \cdot P(x_2^{\xi i} | y = \pm 1) \cdot \dots \cdot P(x_d^{\xi i} | y = \pm 1) \cdot P(y = \pm 1)}{P(x^{\xi i})}$$

univariate gaussian

We move from computing Σ^{-1} which is very computationally expensive to computing product of univariate gaussian

Bayes Decision Rule

- Learning: prior: $p(y)$, class conditional distribution : $p(x|y)$

- The poster probability of a test point

$$\underline{q_i(x)} := P(y = i|x) = \frac{P(x|y)P(y)}{P(x)}$$

- Bayes decision rule:

- If $q_i(x) > q_j(x)$, then $y = i$, otherwise $y = j$

- Alternatively:

- If ratio $l(x) = \frac{P(x|y=i)}{P(x|y=j)} > \frac{P(y=j)}{P(y=i)}$, then $y = i$, otherwise $y = j$

- Or look at the log-likelihood ratio $h(x) = -\ln \frac{q_i(x)}{q_j(x)}$

Generative Model: Naive Bayes

- Use Bayes decision rule for classification

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

- But assume $p(x|y = 1)$ is fully factorized : Dimensions are conditionally independent.

$$p(x|y = 1) = \prod_{i=1}^d p(x_i|y = 1)$$

- Or the variables corresponding to each dimension of the data are independent given the label

“Naive” conditional independence assumption

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} = \frac{P(x, y)}{P(x)}$$

Joint probability model:

$$\begin{aligned} P(x, y_{label=1}) &= P(x_1, \dots, x_d, y_{label=1}) = P(x_1 | x_2, \dots, x_d, y_{label=1}) P(x_2, \dots, x_d, y_{label=1}) \\ &= P(x_1 | x_2, \dots, x_d, y_{label=1}) P(x_2 | x_3, \dots, x_d, y_{label=1}) P(x_3, \dots, x_d, y_{label=1}) \\ &= \dots \\ &= P(x_1 | x_2, \dots, x_d, y_{label=1}) P(x_2 | x_3, \dots, x_d, y_{label=1}) \dots \\ &\quad P(x_{d-1} | x_d, y_{label=1}) P(x_d | y_{label=1}) P(y_{label=1}) \end{aligned}$$

Naïve Bayes assumption: let's rewrite it as:

$$P(x, y_{label=1}) = \underbrace{P(x_1 | y_{label=1})}_{\substack{d \\ \prod_{i=1}^d P(x_i | y_{label=1})}} \underbrace{P(x_2 | y_{label=1})}_{\substack{d \\ \prod_{i=1}^d P(x_i | y_{label=1})}} \dots \underbrace{P(x_d | y_{label=1})}_{\substack{d \\ \prod_{i=1}^d P(x_i | y_{label=1})}} P(y_{label=1}) =$$

Gaussian naïve Bayes
A typical assumption

“Naïve” conditional independence assumption

Real World Example

What do People do in Practice?

$p(x,y)$

- Generative models
 - Model prior and likelihood explicitly
 - “Generative” means able to generate synthetic data points
 - Examples: Naive Bayes, Hidden Markov Models
- Discriminative models
 - Directly estimate the posterior probabilities
 - No need to model underlying prior and likelihood distributions
 - Examples: Logistic Regression, SVM, Neural Networks

Discriminative Models

- Directly estimate decision boundary $h(\mathbf{x}) = -\ln \frac{q_i(\mathbf{x})}{q_j(\mathbf{x})}$ or

posterior distribution $p(y|\mathbf{x})$

- Logistic regression, Neural networks
 - Do not estimate $p(\mathbf{x}|y)$ and $p(y)$
-
- Why discriminative classifier?
 - Avoid difficult density estimation problem  Generative model
 - Empirically achieve better classification results

Recap

$$P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)} = \frac{P(x,y)}{P(x)}$$

$P(y|x)$ → learn this directly

- Generative vs Discriminative

- Generative Example? → Naive Bayes

- Assumptions to calculate posterior → all features are conditionally independent

(Σ^{-1})

$$P(x_1, x_2, \dots, x_d | y) \leftarrow \text{Likelihood}$$

$$\parallel \\ P(x_1 | y) \cdot P(x_2 | y) \cdot \dots \cdot P(x_d | y)$$

Gaussian/Multinomial/Bernoulli Naive Bayes



Bernoulli Naive Bayes

For binary or boolean features.

EXAMPLE



1	0	1
0	1	1
0	1	1
0	1	0
0	0	1



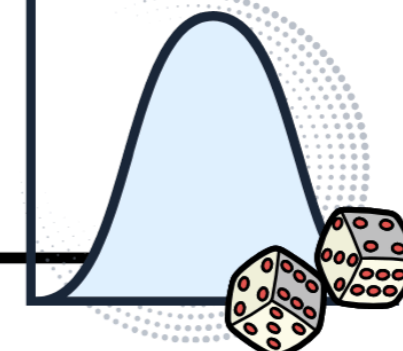
Multinomial Naive Bayes

For discrete features (like word count)

EXAMPLE



4	1	2
3	0	0
3	0	0
1	0	4
0	2	3



Gaussian Naive Bayes

For continuous, real-valued attributes.

EXAMPLE



17.4	56.5	145.2
25.4	71.2	170.4
18.8	70.3	164.5
21.1	51.2	140.5
20.9	81.5	182.2

Outline

- Generative and Discriminative Classification
- The Logistic Regression Model ←
- Understanding the Objective Function ←
- Gradient Descent for Parameter Learning ←
- Multiclass Logistic Regression

2 classes

Gaussian Naïve Bayes

$$P(y = 1|x) = \frac{P(x|y = 1)P(y = 1)}{P(x)} = \frac{P(y = 1) \prod_{i=1}^d P(x_i|y = 1)}{P(x)}$$

fraction of datapoints with y=1

$$\prod_{i=1}^d p(x_i|y = 1, \mu_{1i}, \sigma_{1i})$$
$$= \prod_{i=1}^d \frac{1}{\sqrt{2\pi}\sigma_{1i}} \exp\left(-\frac{1}{2\sigma_{1i}^2} (x_{1i} - \mu_{1i})^2\right)$$

Prior: $p(y = 1) = \pi_1 = \frac{\# \text{ in class 1}}{\text{All data points.}}$

Posterior: $p(y = 1 | x, \mu, \sigma, \pi)$

$$= \frac{\pi_1 \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_{1i}}} \exp\left(-\frac{1}{2\sigma_{1i}^2} (x_i - \mu_{1i})^2\right)}{\sum_{k=1}^2 \pi_k \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_{ki}}} \exp\left(-\frac{1}{2\sigma_{ki}^2} (x_i - \mu_{ki})^2\right)}$$

$\log(abc) = \log a + \log b + \log c$
 $\log\left(\frac{a}{b}\right) = \log a - \log b$

$$\sum_{k=1}^2 \pi_k \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_{ki}}} \exp\left(-\frac{1}{2\sigma_{ki}^2} (x_i - \mu_{ki})^2\right)$$

take $\exp(\ln(u)) = u$ of numerator and denominator

$\log \frac{1}{\sqrt{2\pi}\sigma_{ii}}$
 $\log 1 - \log \sigma_{ii}$

$$= \frac{\exp\left(-\sum_{i=1}^d \left(\frac{1}{2\sigma_{1i}^2} (x_i - \mu_{1i})^2 + \log \sigma_{1i} + C\right) + \log \pi_1\right)}{\sum_{k=1}^2 \exp\left(-\sum_{i=1}^d \left(\frac{1}{2\sigma_{ki}^2} (x_i - \mu_{ki})^2 + \log \sigma_{ki} + C\right) + \log \pi_k\right)}$$

$$P(y=1|x) = \frac{\exp(a)}{\exp(a) + \exp(b)} \quad \text{2 class}$$

$$= \frac{1}{1 + \frac{\exp(b)}{\exp(a)}} = \frac{1}{1 + \exp(b-a)}$$

Assumption: All classes have same variance

$$= \frac{\exp\left(-\sum_{i=1}^d \left(\frac{1}{2\sigma_i^2} (x_i - \mu_{1i})^2 + \log \sigma_i + C\right) + \log \pi_1\right)}{\sum_{k=1}^2 \exp\left(-\sum_{i=1}^d \left(\frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2 + \log \sigma_i + C\right) + \log \pi_k\right)}$$

$\exp(b-a)$ includes: $\frac{(x_i - \mu_{2i})^2 - (x_i - \mu_{1i})^2}{2\sigma_i^2}$

This will lead to no x_i^2 term

$$= \frac{\exp\left(-\sum_{i=1}^d \left(\frac{1}{2\sigma_i^2} (x_i - \mu_{1i})^2 + \log \sigma_i + C\right) + \log \pi_1\right)}{\sum_{k=1}^2 \exp\left(-\sum_{i=1}^d \left(\frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2 + \log \sigma_i + C\right) + \log \pi_k\right)}$$

$$= \frac{1}{1 + \exp\left(-\sum_{i=1}^d \left(\underbrace{x_i \frac{1}{\sigma_i} (\mu_{1i} - \mu_{2i})}_{\text{linear term}} + \underbrace{\frac{1}{\sigma_i^2} (\mu_{1i}^2 - \mu_{2i}^2)}_{\text{bias term}}\right) + \log \frac{\pi_2}{\pi_1}\right)}$$

$1 + \exp\left(\theta_0 + \sum_i \theta_i x_i\right)$

$\sum_i \theta_i x_i$

Sigmoid function θ_0

$$P(y = 1 | x) = \frac{1}{1 + \exp(-(\theta_0 + \sum_i \theta_i x_i))}$$

$$P(y = 1|x) = \frac{1}{1 + \exp(-s)}$$

$$\text{LCF} = \theta_0 + \sum_i \theta_i x_i$$

After all that Gaussian algebra, We started with a generative model...

but we ended up with this.

This is exactly the same form as **logistic regression**.

$$P(y = 1|x) = \frac{1}{1 + \exp\left(-\sum_{i=1}^d \left(x_i \frac{1}{\sigma_i} (\mu_{1i} - \mu_{2i}) + \frac{1}{\sigma_i^2} (\mu_{1i}^2 - \mu_{2i}^2)\right) + \log \frac{\pi_2}{\pi_1}\right)}$$

- Pros:
- ① less parameters
 - ② Make assumptions. Dist. of each class
 - ③ Compute likelihood

Number of parameters:

$2d + 1 \rightarrow d$ mean, d variance, and 1 for prior

$$P(y = 1|x) = \frac{1}{1 + \exp[-(\sum_{i=1}^d (\theta_i x_i) + \theta_0)]} = \frac{1}{1 + \exp(-s)}$$

Number of parameters = $d + 1 \rightarrow \theta_0, \theta_1, \theta_2, \dots, \theta_d$

Why not directly learning $P(y = 1|x)$ or θ parameters?

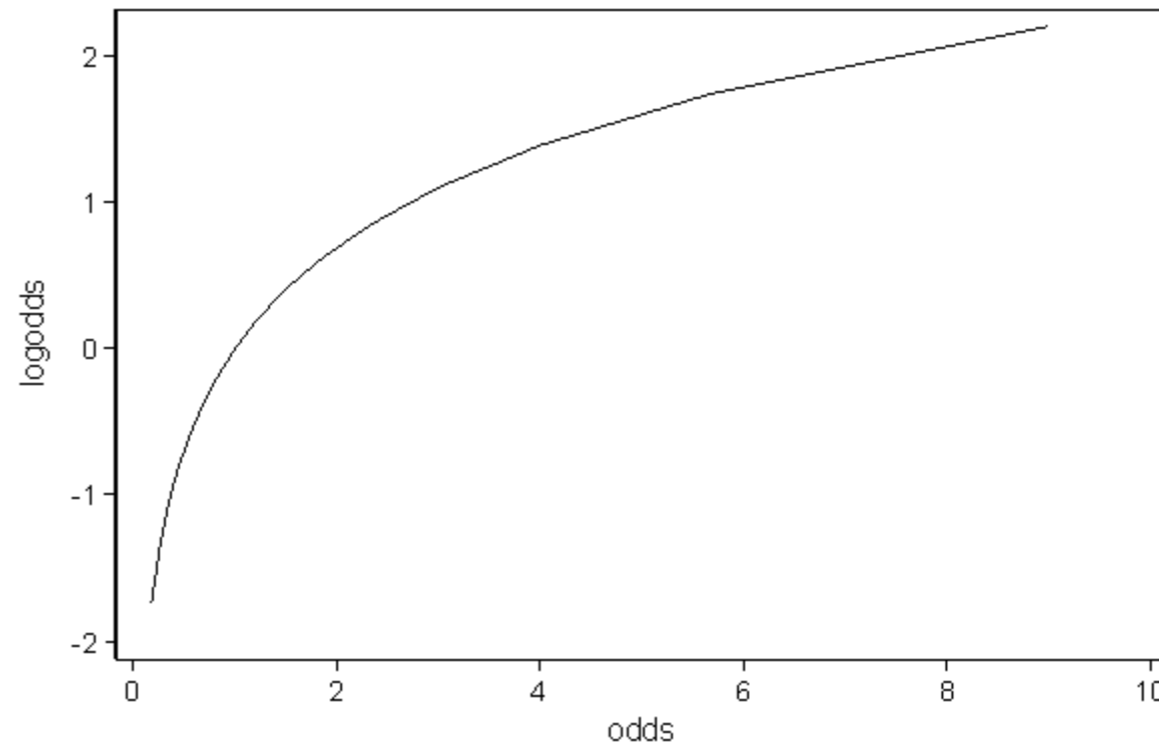
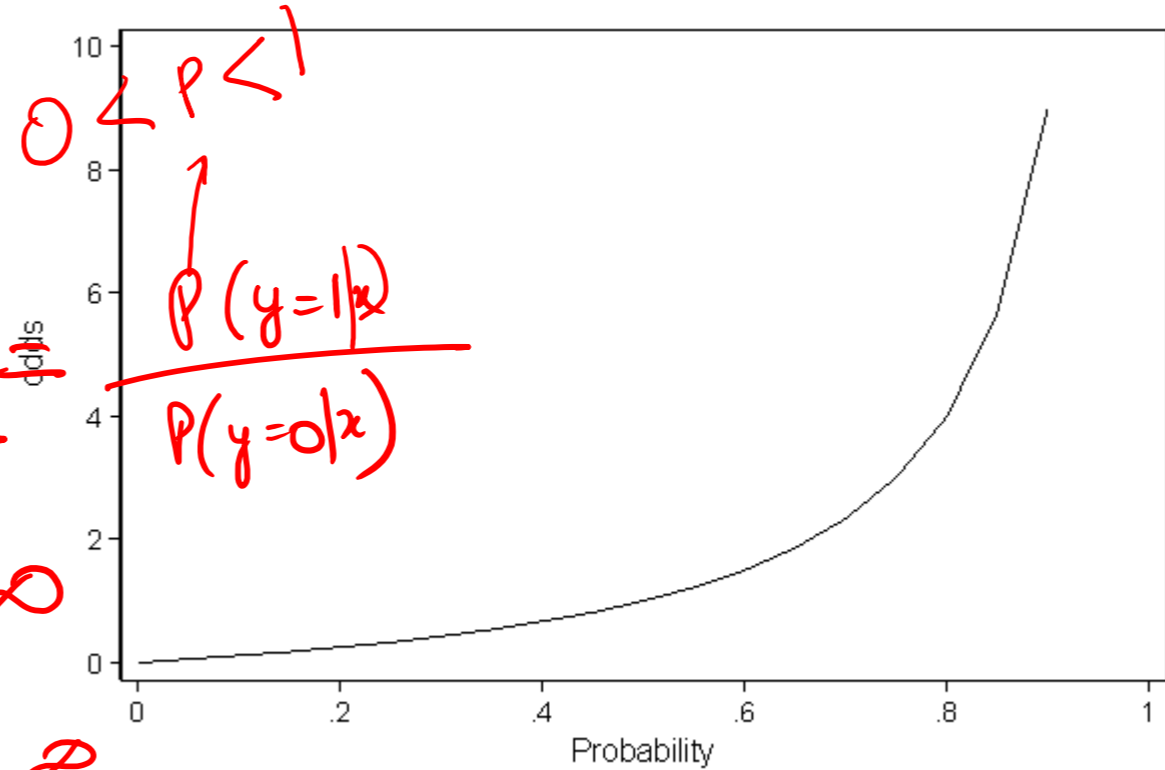
Gaussian Naïve Bayes is a special case/subset of logistic regression

Why $\frac{1}{1+\exp(-x\theta)}$ is a probability?

$\frac{P(y = 1|x)}{1-P(y = 1|x)}$ is called Odds

$0 < \text{odds} < \infty$

$-\infty < \log(\text{odds}) < \infty$



log(odds) vs odds

What could be $x\theta$ domain?

What is logit function?

$$\text{logit}(p) = \log(\text{odds}) = \log\left(\frac{p}{1-p}\right)$$

$-\infty < < \infty$

Assumption

$$\log\left(\frac{p}{1-p}\right) = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d = \sum_{i=0}^d x_i \theta_i = x\theta$$

$$\exp\left(\log\left(\frac{p}{1-p}\right)\right) = \exp(x\theta)$$

Sigmoid function - maps any real-valued $x\theta$ to a probability between 0 and 1.

$$p = \frac{e^{x\theta}}{1 + e^{x\theta}} = \frac{1}{1 + e^{-x\theta}}$$

Interpretation: Generative and Discriminative Thinking

Naive Bayes (Generative):

Models $P(x | y)$



With Gaussian + equal variance:

$$\log \frac{P(y = 1 | x)}{P(y = -1 | x)} = \theta_0 + \sum_i \theta_i x_i$$



Logistic Regression (Discriminative):

$$P(y = 1 | x) = \frac{1}{1 + e^{-(\theta_0 + \sum_i \theta_i x_i)}}$$

Interpretation

$$P(y=1|x) = 0.5 = P(y=-1|x)$$

$$\frac{1}{1 + \exp(-s)} = 0.5 \Rightarrow 1 + \exp(-s) = 2$$

$$\Rightarrow \exp(-s) = 1$$

$$\Rightarrow -s = 0$$

$$\Rightarrow s = 0 \Rightarrow \boxed{w_0 = 0}$$

Linear line

Concept	Gaussian Naïve Bayes	Logistic Regression
Type	Generative	Discriminative
Learns	<u>$P(x,y)$</u> $\rightarrow P(y x)$ and $P(x)$	$P(y x)$
Distribution Assumption	Gaussian features	No distribution assumption
Decision boundary	Linear (after log transformation and equal variance assumption)	Linear
Parameters	Means, variances, and priors $((2d+1))$	Coefficients and bias $((d+1))$

Logistic function for posterior probability

Many equations can give us this shape

Let's use the following function:

$$s = x\theta$$

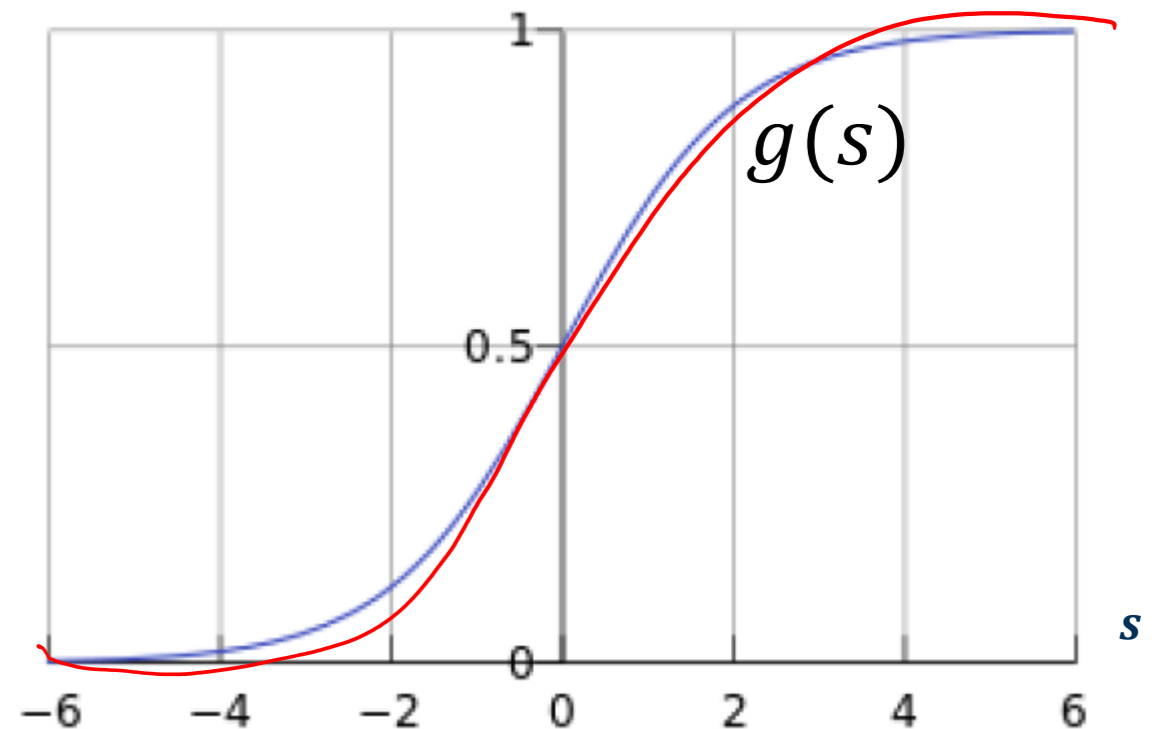
$$\underline{g(s)} = P(y = 1|x) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$

This formula is called sigmoid function

It is easier to use this function for optimization

Is 0.5 threshold cut-off a good choice?

[Learn about ROC and AUC \(False positive rate and True positive rate\) \(Interactive\)](#)



Why Soft classification matters

Flw OR B/w feature selection

Example: Prediction of heart attacks

Input x : cholesterol level, age, weight, finger size, etc.

$g(s)$: probability of heart attack within a certain time

$s = x\theta$ Let's call this risk score

$p=0.3$

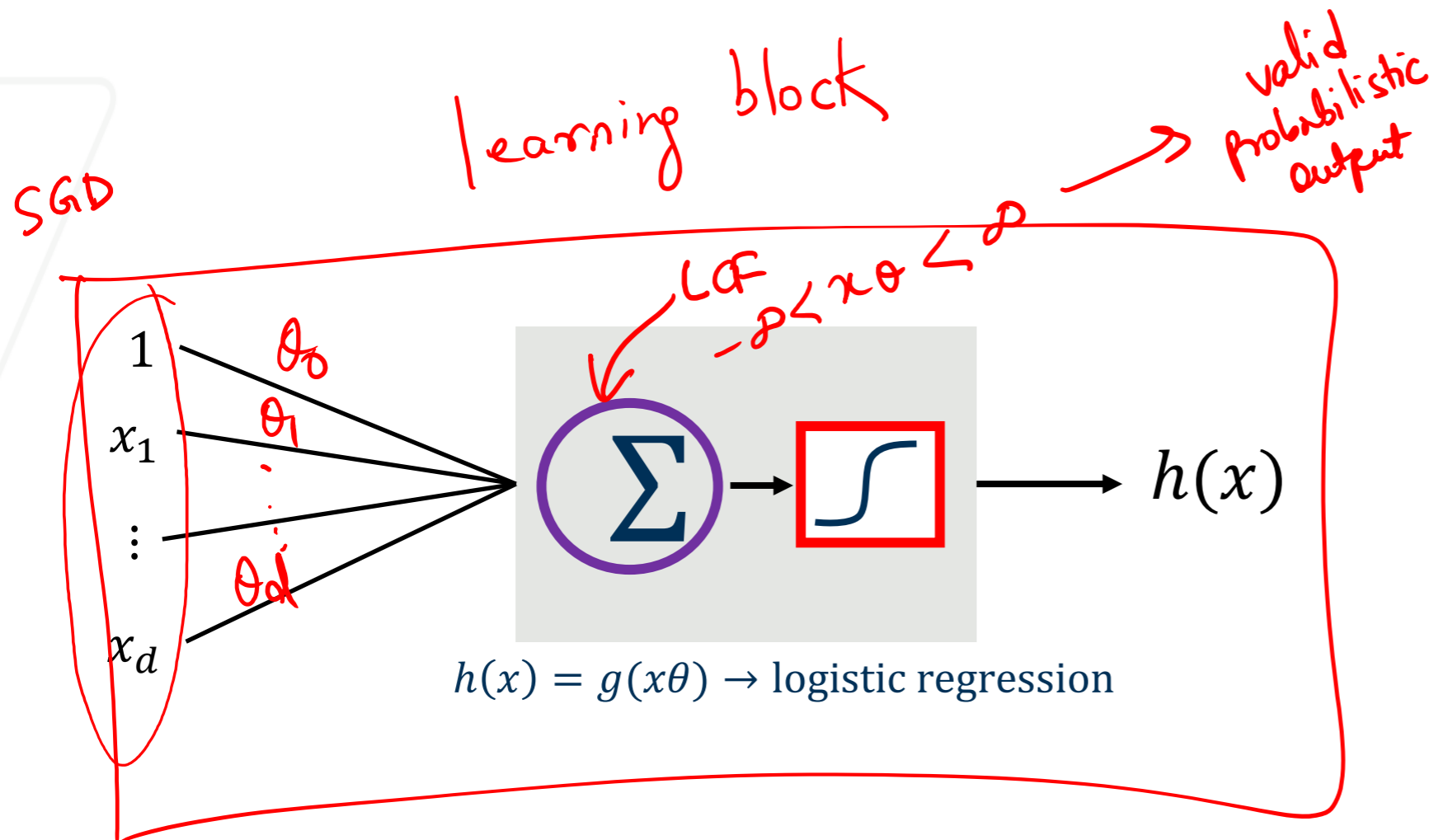
We can't have a hard prediction here

$$h_{\theta}(x) = p(y|x) = \begin{cases} g(s), & y = 1 \\ 1 - g(s), & y = 0 \end{cases}$$
 Using posterior probability directly

Sigmoid Function

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

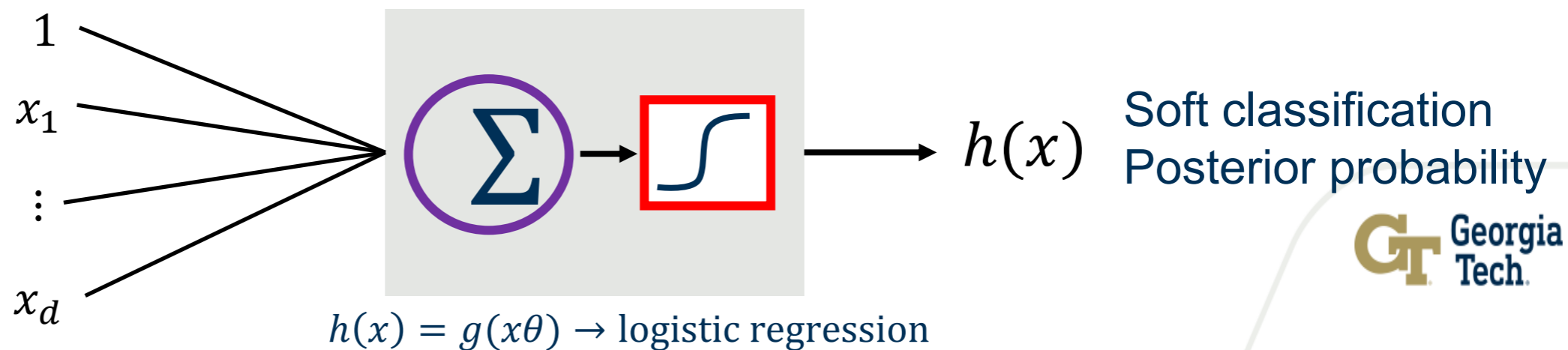
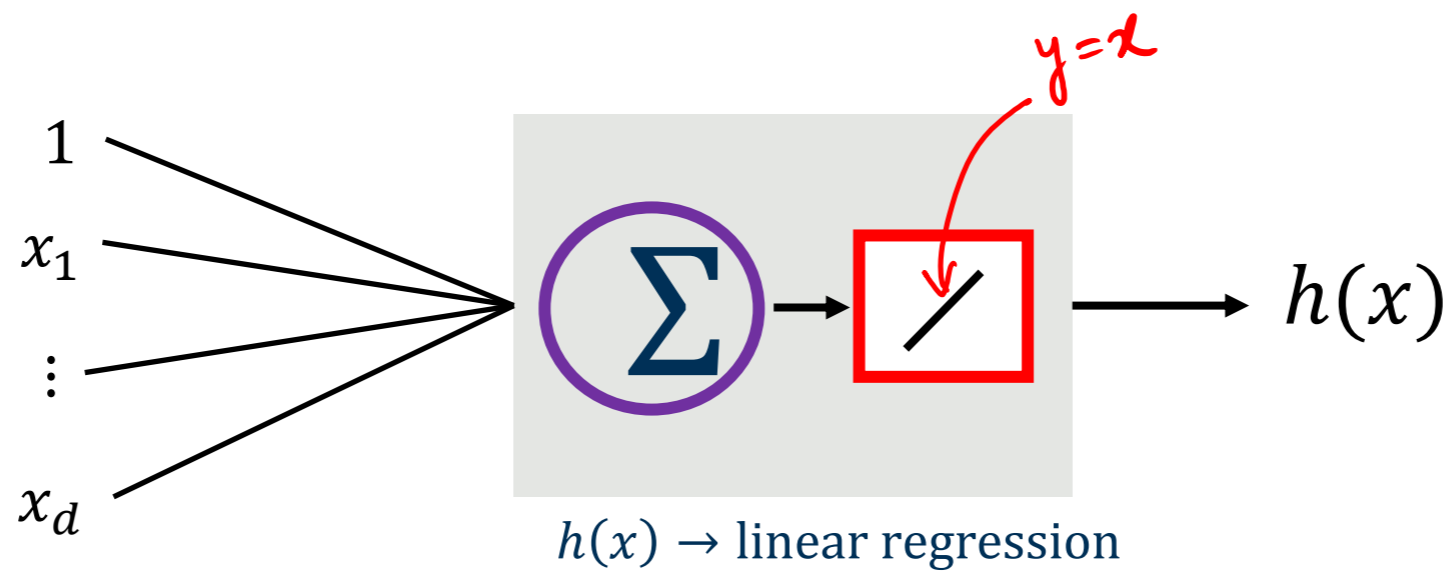
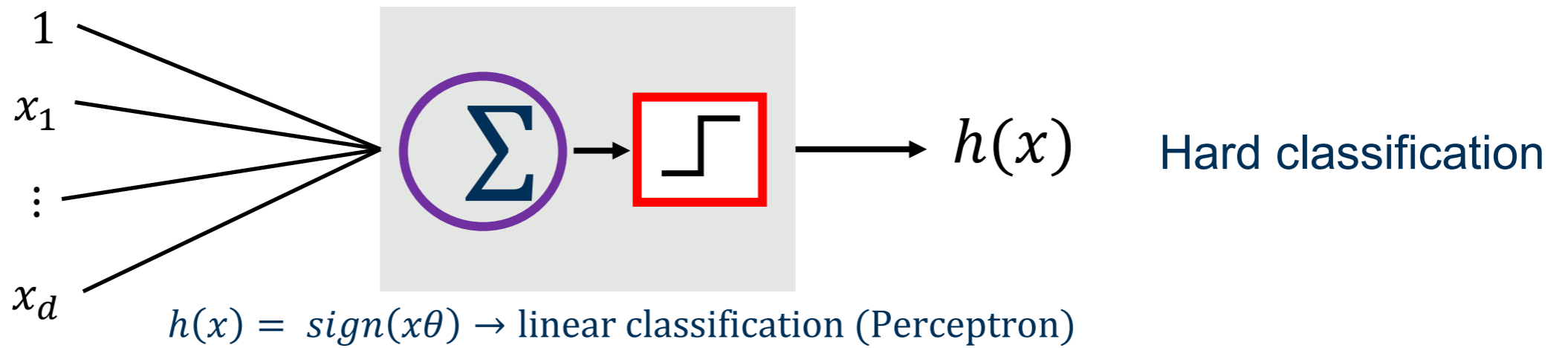
$$s = \sum_{i=0}^d x_i \theta_i = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$



Soft classification
Posterior probability

$$s = \sum_{i=0}^d x_i \theta_i = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

Three linear models



Logistic regression model

$$\underline{p(y|x)} = \begin{cases} \frac{1}{1 + \exp(-x\theta)} & y = 1 \\ 1 - \frac{1}{1 + \exp(-x\theta)} = \frac{\exp(-x\theta)}{1 + \exp(-x\theta)} & y = 0 \end{cases}$$

We need to find θ parameters, what should we use as the objective function?

MLE

We are maximizing the probability of the *observed labels given the inputs*.

We choose θ so that the model assigns high probability to the correct labels.

Logistic regression model

$$p(y|x) = \begin{cases} \frac{1}{1 + \exp(-x\theta)} & y = 1 \\ 1 - \frac{1}{1 + \exp(-x\theta)} = \frac{\exp(-x\theta)}{1 + \exp(-x\theta)} & y = 0 \end{cases}$$

$$p(y|x) = \left(\frac{1}{1 + \exp(-x\theta)} \right)^y \left(\frac{\exp(-x\theta)}{1 + \exp(-x\theta)} \right)^{1-y}$$

We need to find θ parameters, let's set up log-likelihood for n datapoints

iid

$$l(\theta) := \log \prod_{i=1}^n p(y^{(i)} | x^{(i)}, \theta)$$

$$= \sum_i \left(\theta^T (x^{(i)})^T (y^{(i)} - 1) - \log(1 + \exp(-x^{(i)}\theta)) \right)$$

This form is concave, negative of this form is convex

$$\log \prod_i p(y^{(i)} | x^{(i)}, \theta)$$

$$\log \left(\left(\frac{1}{1 + \exp(-x\theta)} \right)^y \left(\frac{\exp(-x\theta)}{1 + \exp(-x\theta)} \right)^{1-y} \right)$$

$$y (\log 1 - \log(1 + \exp(-x\theta))) + (1-y) \log(\exp(-x\theta)) - (1-y) \log(1 + \exp(-x\theta))$$

$$-y \log(1 + \exp(-x\theta)) + \log(\exp(-x\theta)) - y(-x\theta) - \log(1 + \exp(-x\theta))$$

$$+ y \log(1 + \exp(-x\theta))$$

$$\boxed{-x\theta} + xy\theta - \log(1 + \exp(-x\theta))$$

The gradient of $l(\theta)$

$$l(\theta) = \log \prod_{i=1}^n p(y^{i} | x^{i}, \theta)$$
$$= \sum_i \theta^T (x^{i})^T (y^{i} - 1) - \log(1 + \exp(-x^{i} \theta))$$

- Gradient

$$\frac{\partial l(\theta)}{\partial \theta} = \sum_i (x^{i})^T (y^{i} - 1) + (x^{i})^T \frac{\exp(-x^{i} \theta)}{1 + \exp(-x^{i} \theta)}$$

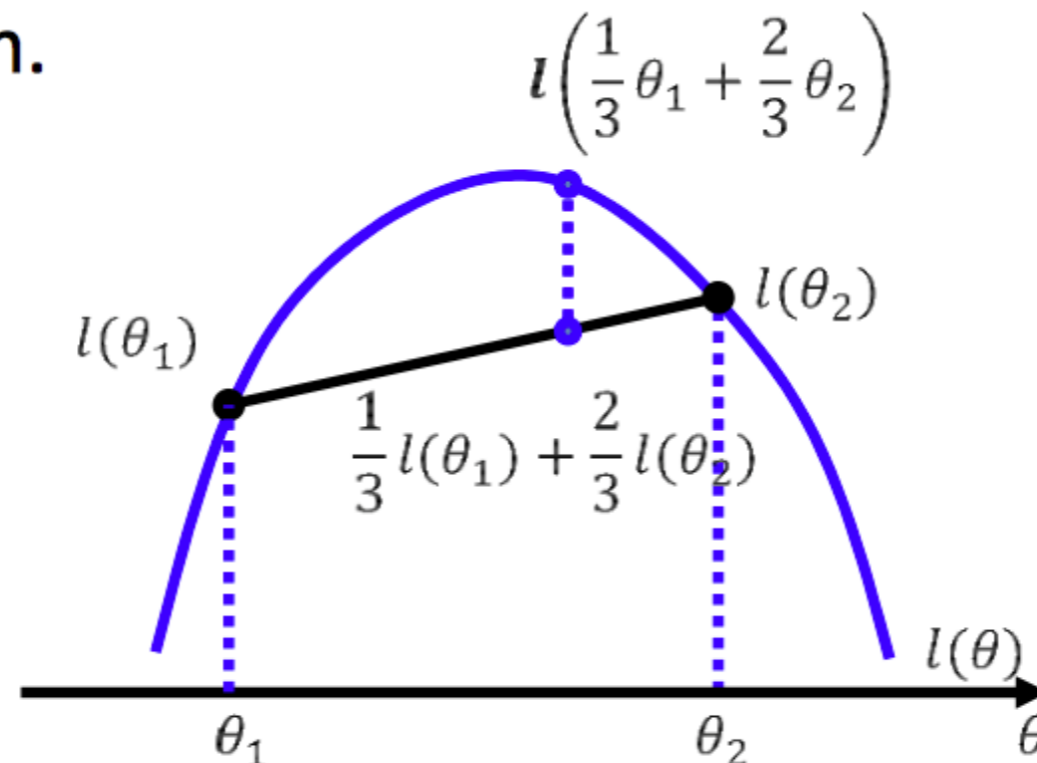
- Setting it to 0 does not lead to closed form solution

The Objective Function

- Find θ , such that the conditional likelihood of the labels is maximized

$$\max_{\theta} l(\theta) := \log \prod_{i=1}^{\bar{n}} p(y^{\{i\}} | x^{\{i\}}, \theta)$$

- Good news: $l(\theta)$ is concave function of θ , and there is a single global optimum.



- Bad news: no closed form solution (resort to numerical method)

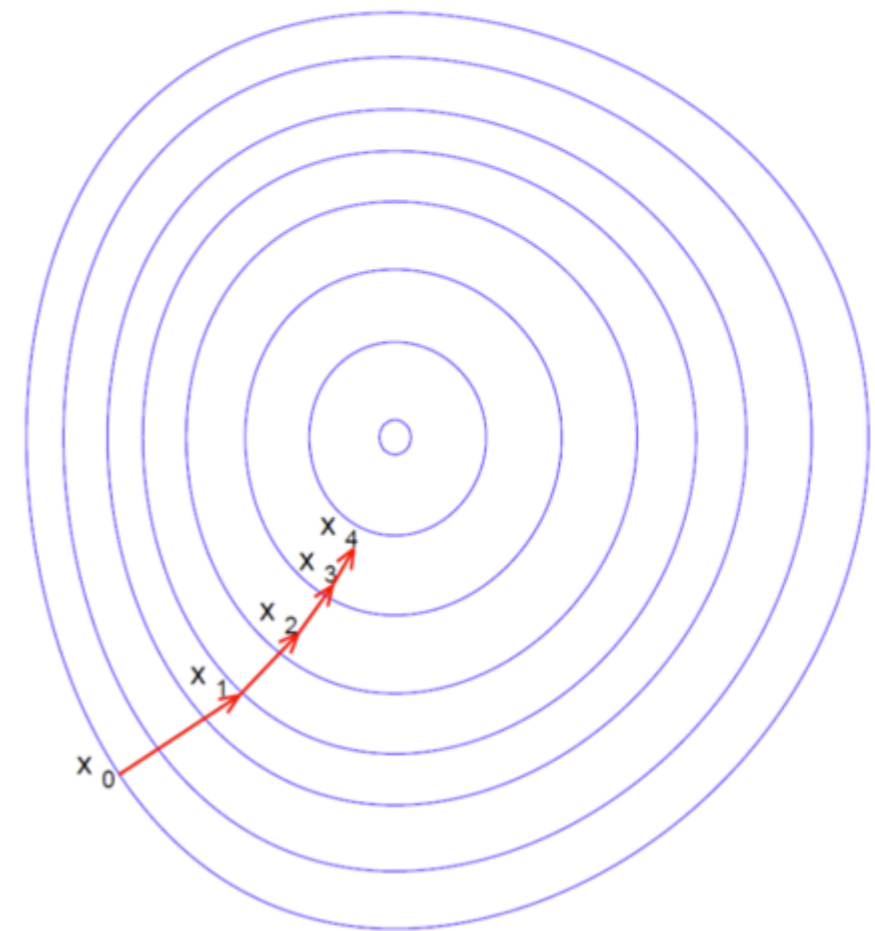
Gradient Descent

- One way to solve an *unconstrained* optimization problem is gradient descent
- Given an initial guess, we *iteratively* refine the guess by taking the direction of the negative gradient
- Think about going down a hill by taking the steepest direction at each step

- Update rule

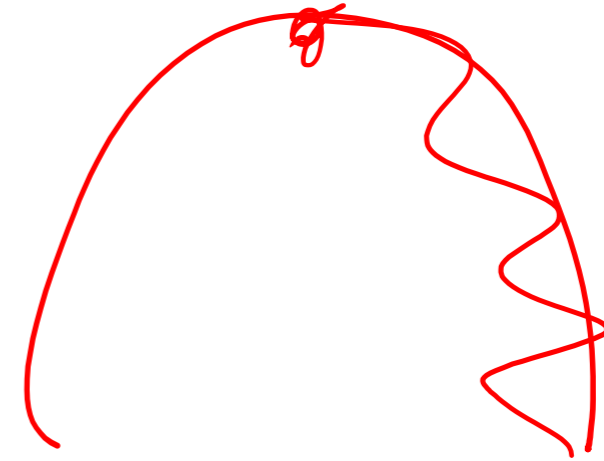
$$x_{k+1} = x_k - \gamma_k \nabla f(x_k)$$

γ_k is called the step size or learning rate



Gradient Ascent(concave)/Descent(convex) algorithm

- Initialize parameter θ^0
- Do

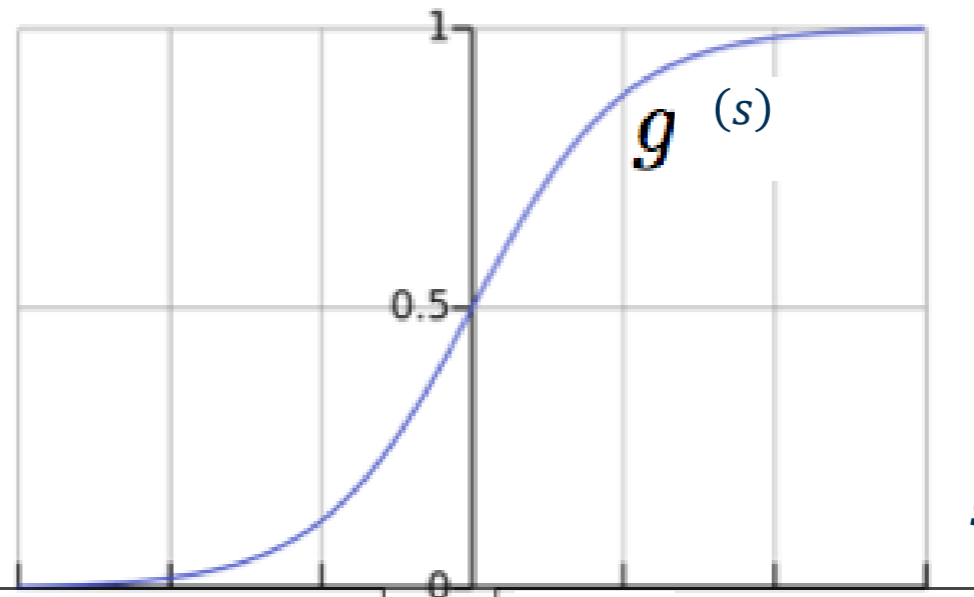


$$\theta^{t+1} \leftarrow \theta^t + \eta \sum_i (x^{i})^T (y^{i} - 1) + (x^{i})^T \frac{\exp(-x^{i}\theta)}{1 + \exp(-x^{i}\theta)}$$

- While the $\|\theta^{t+1} - \theta^t\| > \epsilon$

Logistic Regression

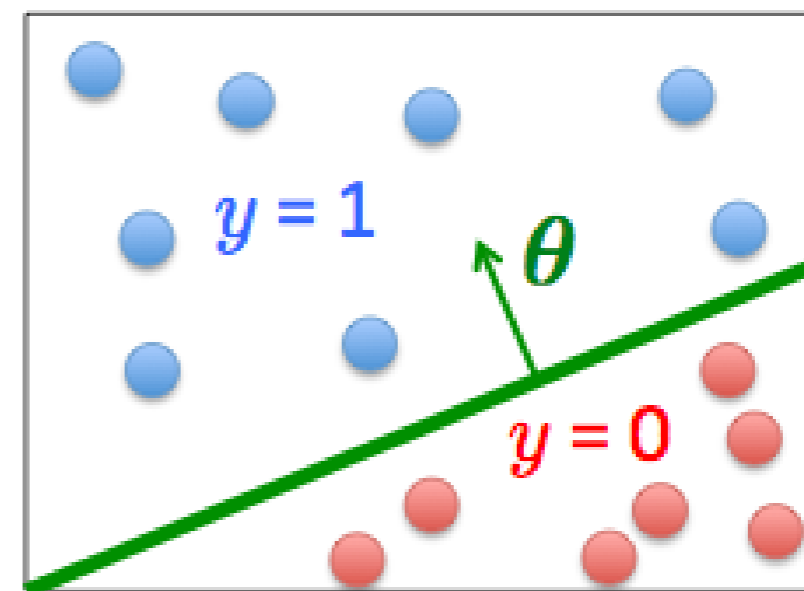
$$g(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$
$$s = x\theta$$



$x\theta$ should be large negative values for negative instances

$x\theta$ should be large positive values for positive instances

- Assume a threshold and...
 - Predict $y = 1$ if $g(s) \geq 0.5$
 - Predict $y = 0$ if $g(s) < 0.5$

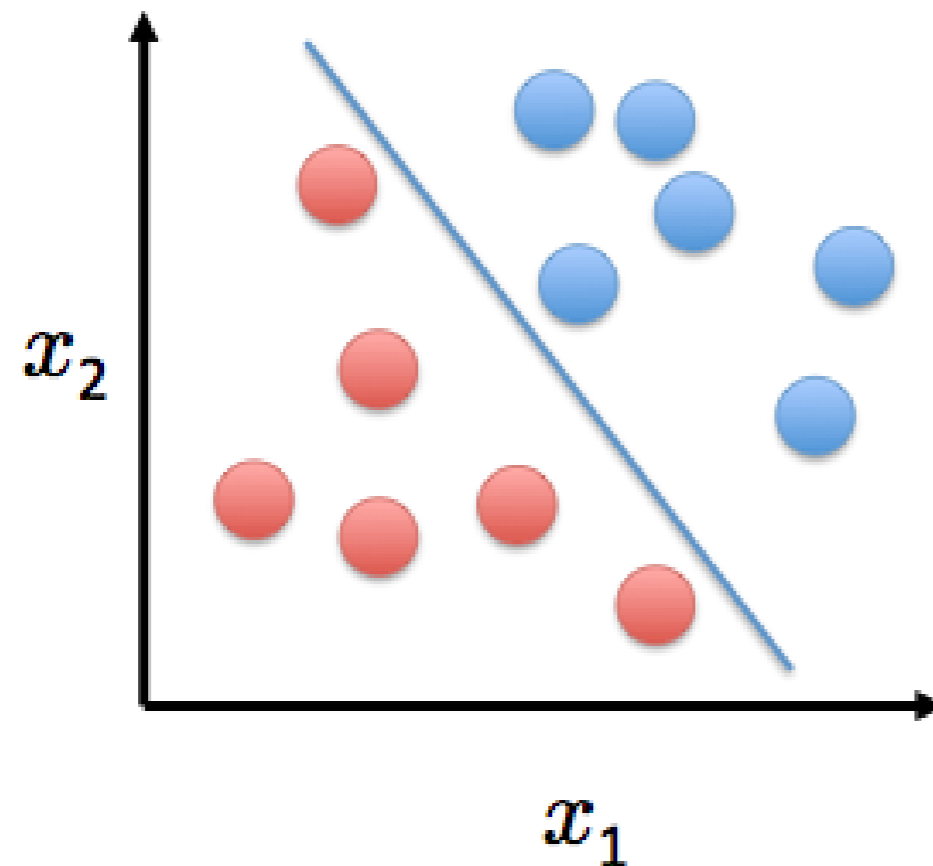


Outline

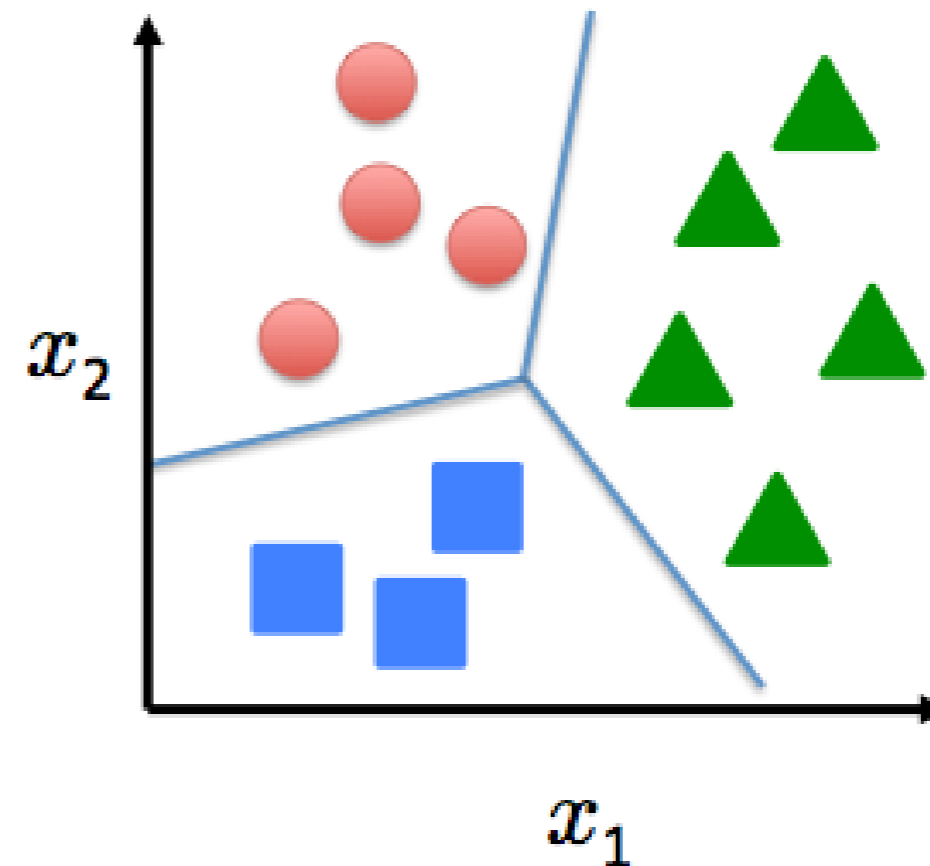
- Generative and Discriminative Classification
- The Logistic Regression Model
- Understanding the Objective Function
- Gradient Descent for Parameter Learning
- Multiclass Logistic Regression ←

Multiclass Logistic Regression

Binary classification:



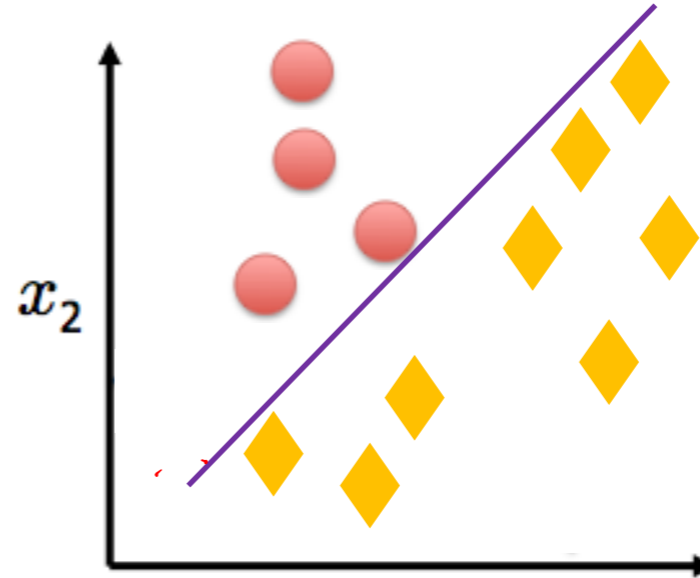
Multi-class classification:



Disease diagnosis: healthy / cold / flu / pneumonia

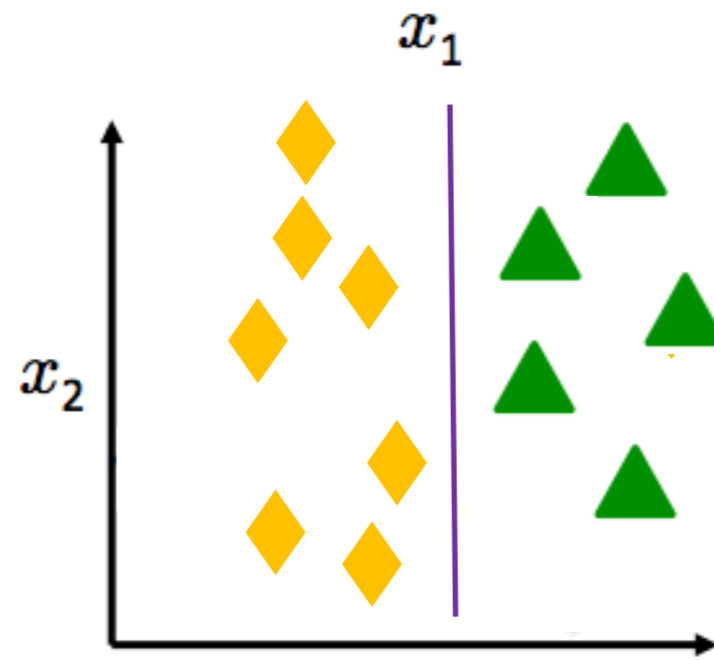
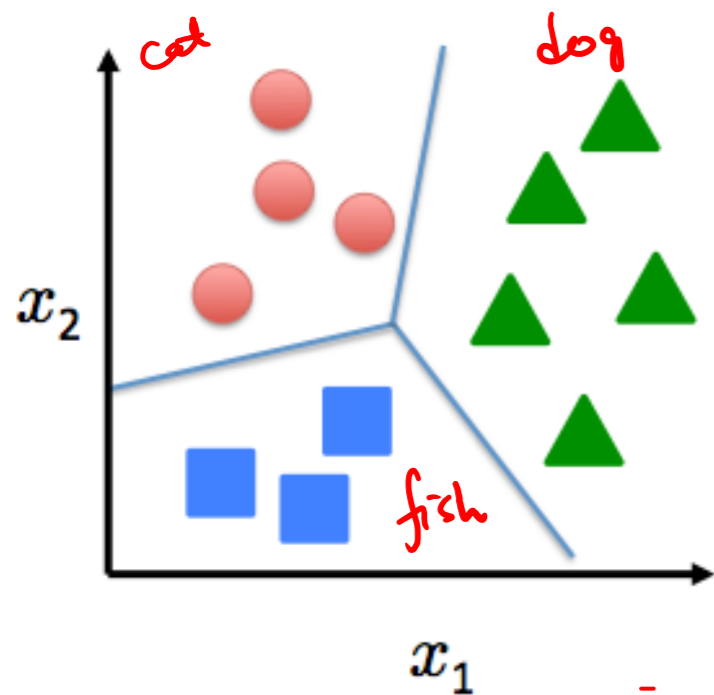
Object classification: desk / chair / monitor / bookcase

One-vs-all (one-vs-rest)

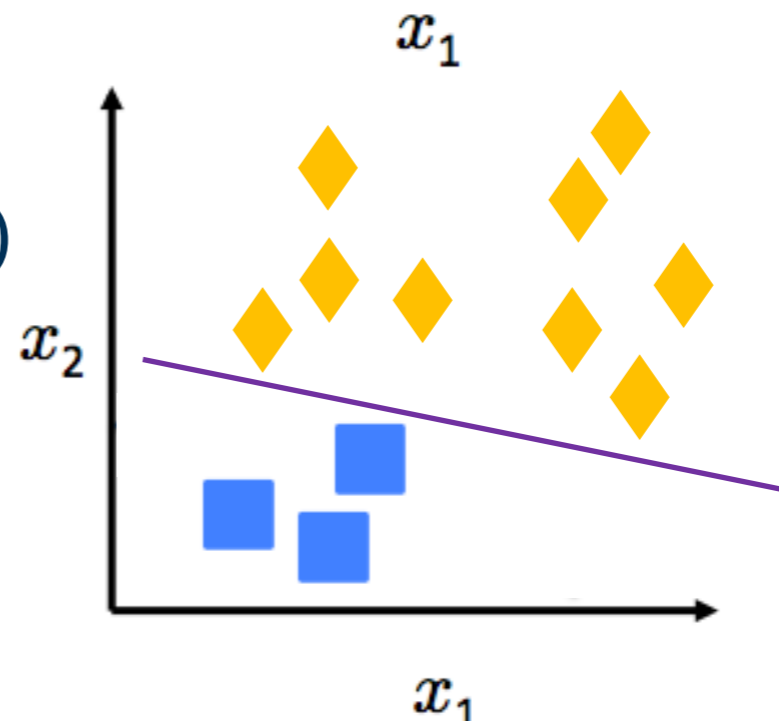


$h_{\theta}^1(x)$
 $\{ 0.7 \quad 0.3 \}$
 Cat not cat

Multi-class classification:



$h_{\theta}^2(x)$
 $\{ 0.6 \quad 0.4 \}$
 dog not dog



$h_{\theta}^3(x)$
 $\{ 0.8 \quad 0.2 \}$
 fish not fish

$$h_{\theta}^{(m)}(x) = p(y = 1 | x, \theta) \quad (m = 1, 2, 3)$$

One-vs-all (one-vs-rest)

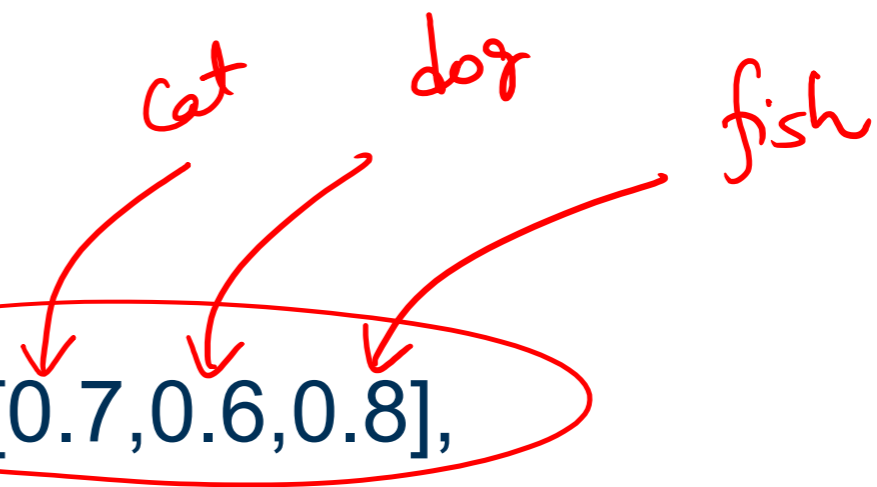
Train a logistic regression $h_{\theta}^{(m)}(x)$ for each class m

To predict the label of a new input x , pick class m that maximizes:

$$\arg \max_i h_{\theta}^{(m)}(x)$$

Key properties:

- Output is not a true distribution e.g. $[0.7, 0.6, 0.8]$, because Models are independent
- Need to train m models



Approach 2: Using Cross Entropy + Softmax

Idea:

Train one model

All classes compete together

Output:

$$p_m = \frac{\exp((x\theta)_m)}{\sum_j \exp((x\theta)_j)}$$

Key properties:

Outputs:

[0.6, 0.3, 0.1] (sums to 1)

True probability distribution

Approach 2: Using Cross Entropy + Softmax

$$L(\theta) = - \sum_{i=1}^N y_a^{\{i\}} * \log(y_p^{\{i\}})$$

Handwritten annotations in red: $y_a^{\{i\}}$ is circled and has arrows pointing to $[0, 0]$ and $[0.7, 0.3]$. $\log(y_p^{\{i\}})$ is circled and has arrows pointing to $[0.7, 0.3]$ and $[0.2]$. The entire equation is circled.

Objective function - Make the probability of the correct class high

Sigmoid (binary case)

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

Outputs: One probability p . Other class is $1 - p$

Softmax (multiclass case)

$$p_m = \frac{\exp(s_m)}{\sum_j \exp(s_j)} = \frac{\exp(s_0)}{\exp(s_0) + \exp(s_1)}$$

Handwritten red annotations: $= \frac{1}{1 + \exp(s_1 - s_0)}$ with an arrow pointing from the right side of the equation to the sigmoid function above.

Outputs: One probability per class, All sum to 1

$[cat, dog, fish]$ $d=2$ $[]$

$x^{(i)} = [2, 1]$

Using Cross Entropy + Softmax

Setup (Dimensions)

$x \in \mathbb{R}^d$ input features

$\theta \in \mathbb{R}^{d \times M}$ parameters

$s \in \mathbb{R}^M$ scores = $x \cdot \theta$ for each class

$y_a \in \mathbb{R}^M$ one-hot label

$y_p \in \mathbb{R}^M$ predicted probabilities

Step 1: Compute scores (LCF)

$$s = x^T \theta \ (\in \mathbb{R}^M)$$

$\theta_{2 \times 3} = \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 1 \end{bmatrix}$

$(2 \ 1) \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 1 & -1 \end{bmatrix}$

cat dog fish

Step 2: Convert to probabilities (Softmax)

$$p_m = \frac{\exp(s_m)}{\sum_j \exp(s_j)} \text{ for } m = 1, \dots, M$$

$$y_p \in \mathbb{R}^M, \sum_m p_m = 1$$

$$p_{cat} = \frac{\exp(s_{cat})}{\exp(s_{cat}) + \exp(s_{dog}) + \exp(s_{fish})}$$

$$= \frac{\exp(2)}{\exp(2) + \exp(1) + \exp(-1)}$$
$$= 0.7$$

$$p_{dog} = 0.25$$

$$p_{fish} = 0.05$$

Using Cross Entropy + Softmax

$$y_a = [1 \ 0 \ 0]$$
$$\hat{y}_p = p_m = [0.7 \ 0.25 \ 0.05]$$

Step 3: Compute loss (Cross-Entropy)

$$L = - \sum_{m=1}^M y_a^{(m)} \log(p_m)$$

• For one datapoint: $L = -\log(p_{\text{correct class}})$

$$L = -\log(0.7)$$

Step 4: Gradient

$$\nabla_{\theta} L = x(y_p - y_a)^T \quad (\in \mathbb{R}^{d \times M})$$

2×1 3×1 3×1

$\nabla_{\theta} L$ is of size 2×3

$$\Theta_{2 \times 3} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \cdot [1 - 0.7 \quad 0 - 0.25 \quad 0 - 0.05]$$

Step 5: SGD Update

$$\theta \leftarrow \theta - \alpha (x(y_p - y_a)^T)$$

2×3 2×3 2×3

hyp.

$$\Theta_{\text{new}} = \begin{bmatrix} 1.5 & -0.5 & -0.01 \\ 0.5 & 0.5 & 0.2 \end{bmatrix}$$

Prediction (after training)

$$\hat{y} = \arg \max_m p_m$$

Take-Home Messages

- Generative and Discriminative Classification
- The Logistic Regression Model
- Understanding the Objective Function – Log Likelihood with sigmoid
- Gradient Descent for Parameter Learning
- Multiclass Logistic Regression – Using CE as objective function

Quick Knowledge Check

- Which of the following best describes a generative model? A. Directly models $P(y|x)$ B. Models $P(x|y)$ and $P(y)$, then uses Bayes' rule for $P(y|x)$ C. Uses cross-entropy to learn the decision boundary D. Requires no assumptions about data distribution
- Logistic Regression is an example of a: A. Generative model that reconstructs data distribution B. Hybrid model using both priors and likelihoods C. Discriminative model that directly estimates $P(y|x)$ D. Model that approximates $P(x|y)$ with Gaussian likelihoods
- What does the logit function represent? A. The derivative of the sigmoid B. The log of the odds ratio $\log \frac{P(y=1|x)}{1 - P(y=1|x)}$ C. The posterior probability itself D. The exponential of the decision boundary
- In binary logistic regression, the model uses the sigmoid $g(x) = \frac{1}{1 + e^{-x\theta}}$ to predict class probabilities. The objective function typically: A. Maximizes the joint likelihood $P(x,y)$ B. Minimizes squared error between predictions and labels C. Minimizes the negative log-likelihood based on the sigmoid output D. Minimizes the variance of the Gaussian features

$$\text{Odds} = \frac{P(y=1)}{P(y=0)} = \text{Logit}$$

Quick Knowledge Check

- In multiclass logistic regression using softmax, the objective function: A. Computes separate sigmoids for each class independently B. Minimizes the mean squared error across classes C. Maximizes the total log-likelihood using cross-entropy over all classes
- Gradient descent in logistic regression is used to: A. Optimize parameters by minimizing or maximizing the objective function B. Update weights in the direction of the gradient of log-likelihood C. Compute priors and posteriors explicitly D. Eliminate class imbalance by reweighting samples
- True or False: Gaussian Naïve Bayes and Logistic Regression yield identical decision boundaries under equal variance and Gaussian feature assumptions. A. True B. False