

Naïve Bayes and Logistic Regression

Mahdi Roozbahani
Georgia Tech



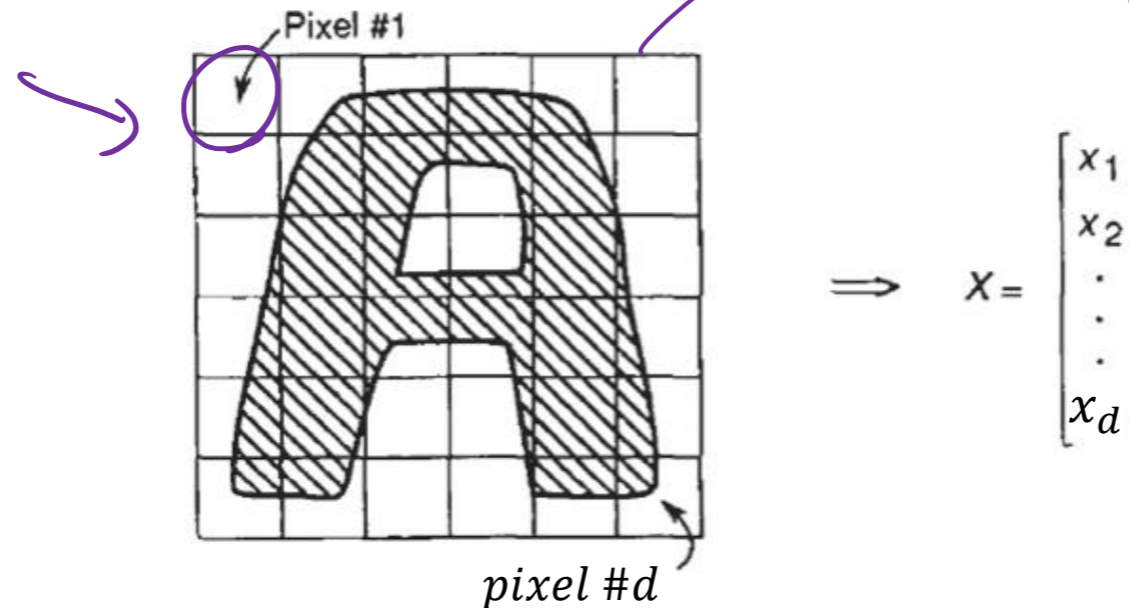
**THE BEST WAY TO
EXPLAIN OVERFITTING**

Outline

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

Classification

- Represent the data

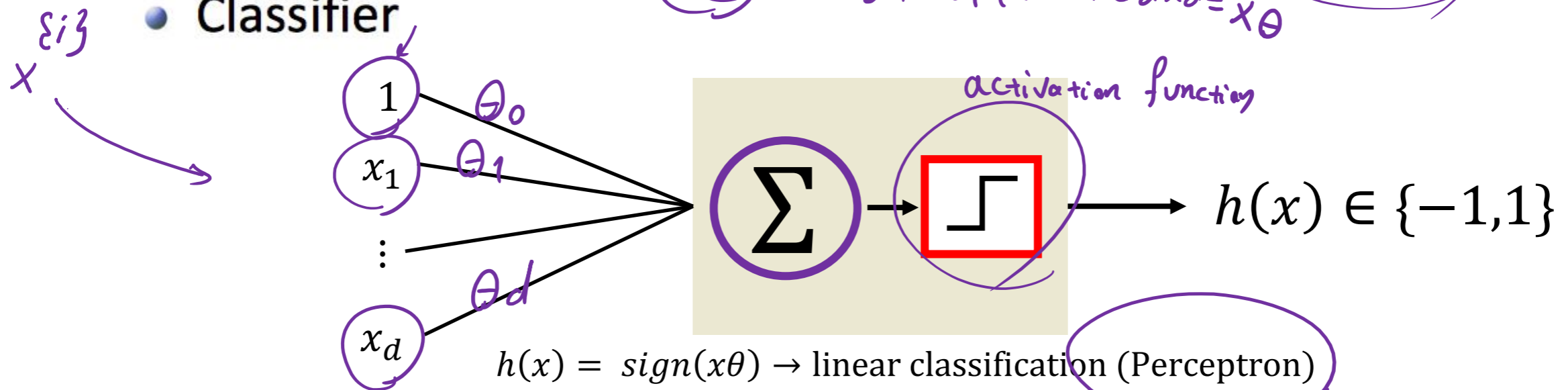


- A label is provided for each data point, eg., $y \in \{-1, +1\}$

- Classifier

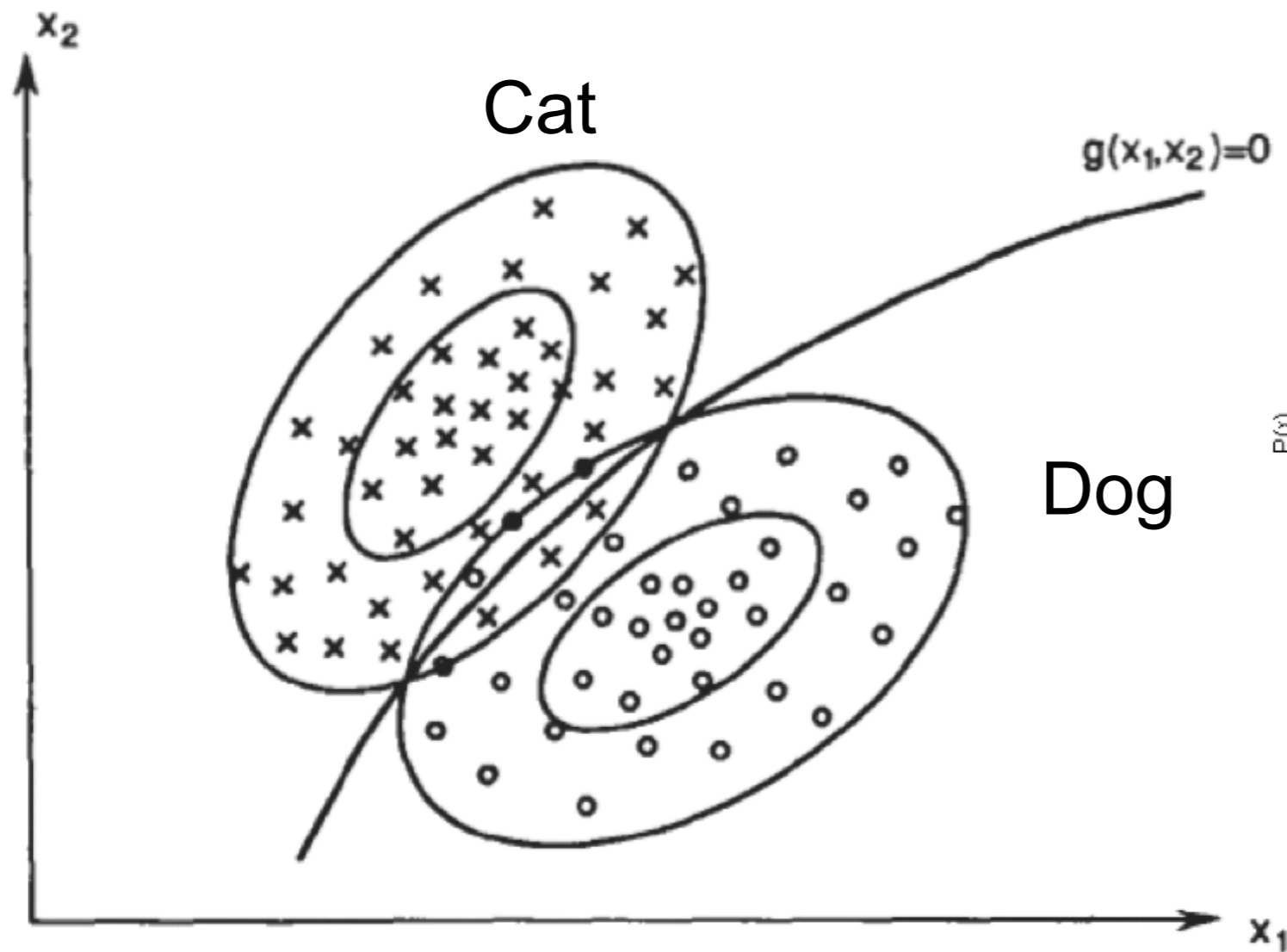
$\Sigma = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d = x\theta$

$\begin{matrix} 0 & 1 \end{matrix}$

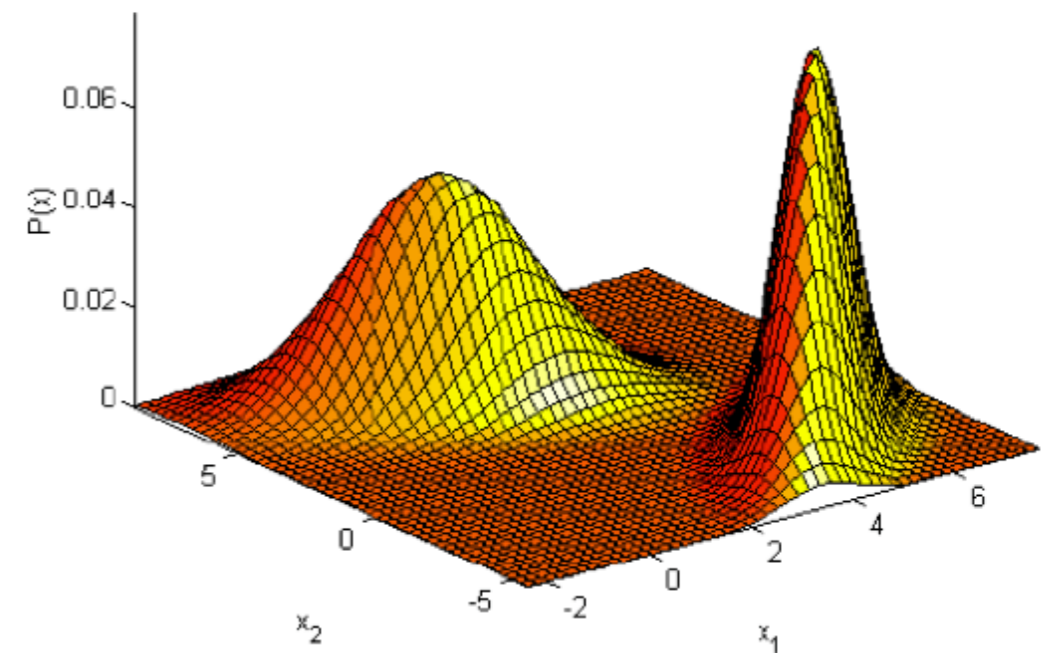


Decision Making: Dividing the Feature Space

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

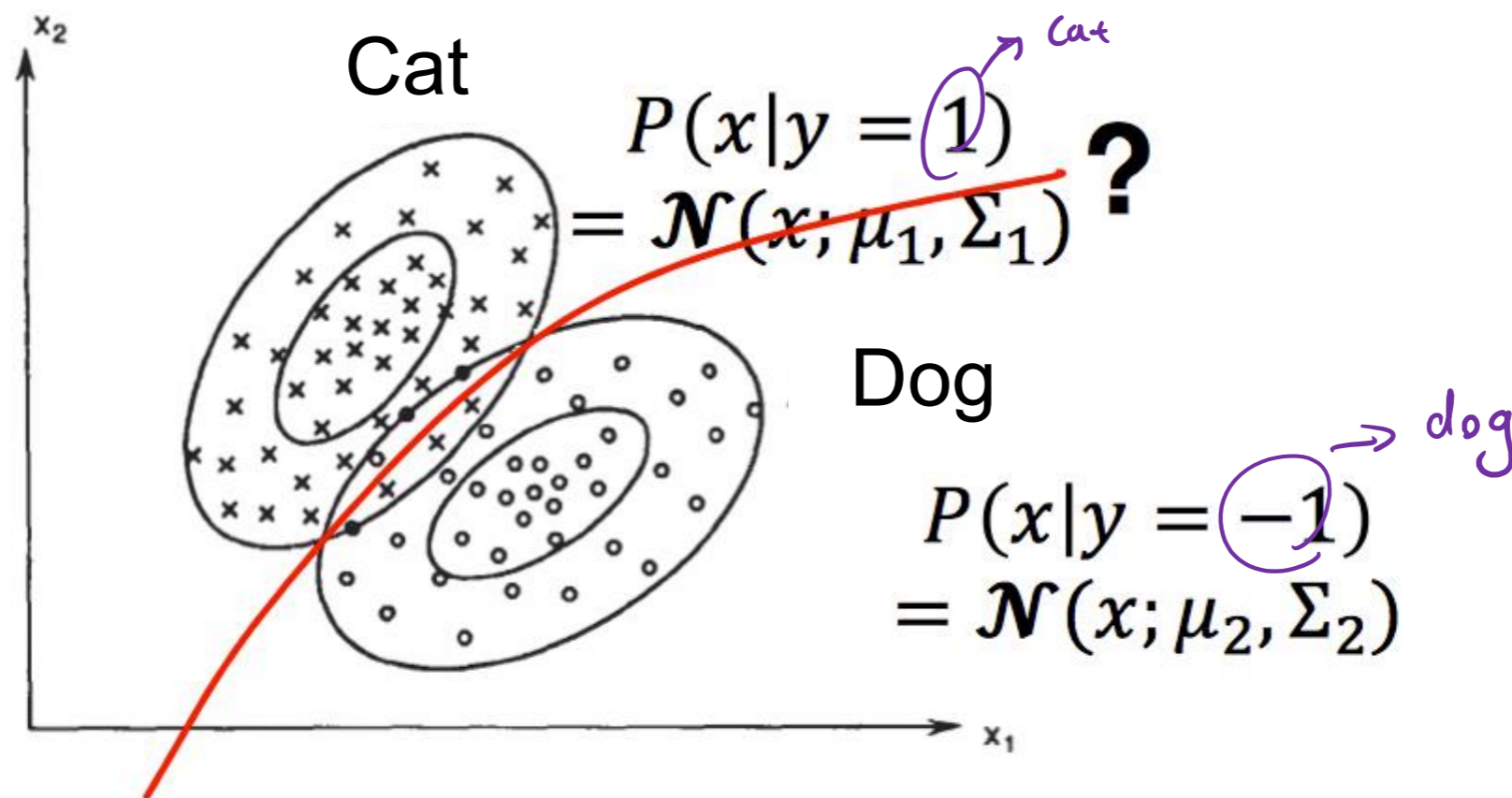


$$X = \begin{bmatrix} h & w \end{bmatrix} \quad Y = \begin{bmatrix} \text{cat} \\ i \\ \text{dog} \\ i \end{bmatrix}$$



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

The diagram shows the equation for Bayes' theorem with four red arrows pointing to its components: 'likelihood' points to $P(x|y)$, 'Prior' points to $P(y)$, 'posterior' points to $P(y|x)$, and 'normalization constant' points to $P(x)$.

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

Prior: $P(y)$

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

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

Bayes Decision Rule

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

- The poster probability of a test point

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

- 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 $h(x) = -\ln \frac{q_i(x)}{q_j(x)}$

$$X = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix} \quad Y = \begin{bmatrix} \text{cat} \\ \text{cat} \\ \text{cat} \\ \text{dog} \\ \text{dog} \\ \text{dog} \end{bmatrix}$$

$$X^{\{\pm\}} = \begin{bmatrix} h & w \\ 13 & 14 \end{bmatrix} \quad Y = ?$$

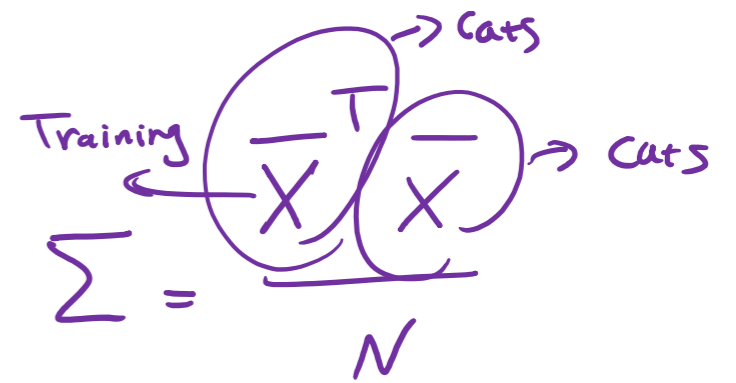
$$Y = +1 \rightarrow \text{cat} \quad Y = -1 \rightarrow \text{dog}$$

$$M = [M_h \quad M_w]$$

$$P(Y = +1 | X^{\{\pm\}}) = \frac{P(Y = +1, X^{\{\pm\}})}{P(X^{\{\pm\}})} = \frac{\overbrace{P(X^{\{\pm\}} | Y = +1)}^{\text{likelihood}} \underbrace{P(Y = +1)}_{\text{prior}}}{P(X^{\{\pm\}}) = \sum_y P(X^{\{\pm\}}, y) = \underbrace{P(X^{\{\pm\}}, y = +1)}_a + \underbrace{P(X^{\{\pm\}}, y = -1)}_b}$$

$$P(Y = -1 | X^{\{\pm\}}) = \frac{b}{a+b}$$

$$P(Y=+1 | x^{\{+\}}) = \frac{N(x | \mu, \Sigma) P(Y=+1)}{P(x^{\{+\}})}$$



Cats

$$N(x^{\{+\}} | \mu, \Sigma) = \frac{1}{\sqrt{2\pi \Sigma^{d/2}}} \exp \left[(x^{\{+\}} - \mu) \Sigma^{-1} (x^{\{+\}} - \mu)^T \right]$$

$$P(Y=-1 | x^{\{+\}}) = ?$$

all features are conditionally independent

Naïve Bayes

$$\text{Cov} = \begin{bmatrix} \sigma_n^2 & 0 \\ 0 & \sigma_w^2 \end{bmatrix} \Sigma^{-1}$$

$$X^{\{+\}} = [x_1, \dots, x_d]$$

$$P(Y=+1 | X^{\{+\}})$$

$$= \frac{P(x_1^{\{+\}}, x_2^{\{+\}}, \dots, x_d^{\{+\}} | Y=+1) P(Y=+1)}{P(X^{\{+\}})}$$

$$= \frac{P(x_1^{\{+\}} | Y=+1) P(x_2^{\{+\}} | Y=+1) \dots P(x_d^{\{+\}} | Y=+1) P(Y=+1)}{P(X^{\{+\}})}$$

What do People do in Practice?

- 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

$$P(Y = +1 | X^{\{+3\}}) \Rightarrow \text{directly}$$

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

“Naïve” conditional independence assumption

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

Joint probability model:

$$P(x, y_{label=1}) = P(x_1, \dots, x_d, y_{label=1}) = P(x_1 | x_2, \dots, x_d, y_{label=1}) P(x_2, \dots, x_d, y_{label=1})$$

$$= P(x_1 | x_2, \dots, x_d, y_{label=1}) P(x_2 | x_3, \dots, x_d, y_{label=1}) P(x_3, \dots, x_d, y_{label=1})$$
$$= \dots$$

$$= P(x_1 | x_2, \dots, x_d, y_{label=1}) P(x_2 | x_3, \dots, x_d, y_{label=1}) \dots P(x_{d-1} | x_d, y_{label=1}) P(x_d | y_{label=1}) P(y_{label=1})$$

Naïve Bayes assumption: let's rewrite it as:


$$P(x, y_{label=1}) = P(x_1 | y_{label=1}) P(x_2 | y_{label=1}) \dots P(x_d | y_{label=1}) P(y_{label=1}) =$$

$$P(y_{label=1}) \prod_{i=1}^d P(x_i | y_{label=1})$$

Gaussian naïve Bayes
A typical assumption

[Example](#)

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

Outline

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

Gaussian Naïve Bayes GNB

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

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

Prior: $p(y = 1) = \pi_1$

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

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

get $\exp(\ln(u))$ of numerator and denominator

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

$$\frac{a}{a+b} \frac{a/a}{a/a + b/a} = \frac{1}{1+b/a}$$

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

$$P(y=+1|x) = \frac{1}{1 + \exp\left(-\sum_{i=1}^d \left(\underbrace{x_i \frac{1}{\sigma_i} (\mu_{1i} - \mu_{2i})}_{\theta_i} + \underbrace{\frac{1}{\sigma_i^2} (\mu_{1i}^2 - \mu_{2i}^2)}_{\theta_0}\right) + \log \frac{\pi_2}{\pi_1}\right)}$$

$$= \frac{1}{1 + \exp\left(-\left[\sum_{i=1}^d \theta_i x_i + \theta_0\right]\right)} = \frac{1}{1 + \exp\left[-\underbrace{(\theta_0 + \theta_1 x_1 + \dots + \theta_d x_d)}_{LCF}\right]}$$

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

Number of parameters:

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

Summation \downarrow

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

$S = LCF = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d = x\theta$

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

Sigmoid $\left[\begin{matrix} 1 \\ 0 \end{matrix} \right]$

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

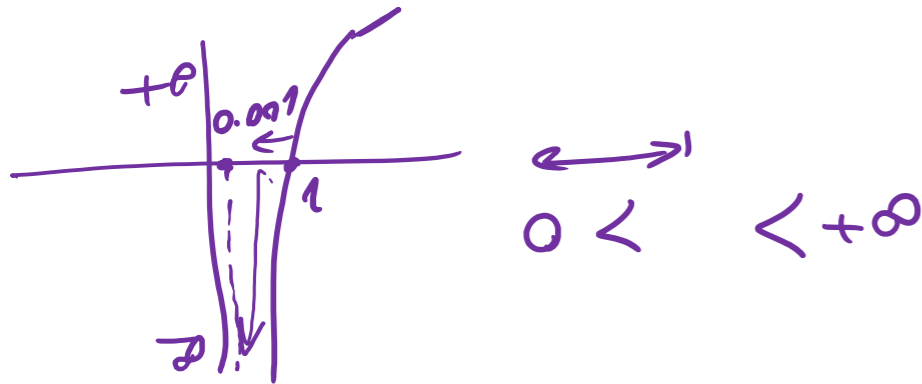
Gaussian Naïve Bayes is a subset of logistic regression

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

$$X\theta = LCF = S \in \mathbb{R}$$

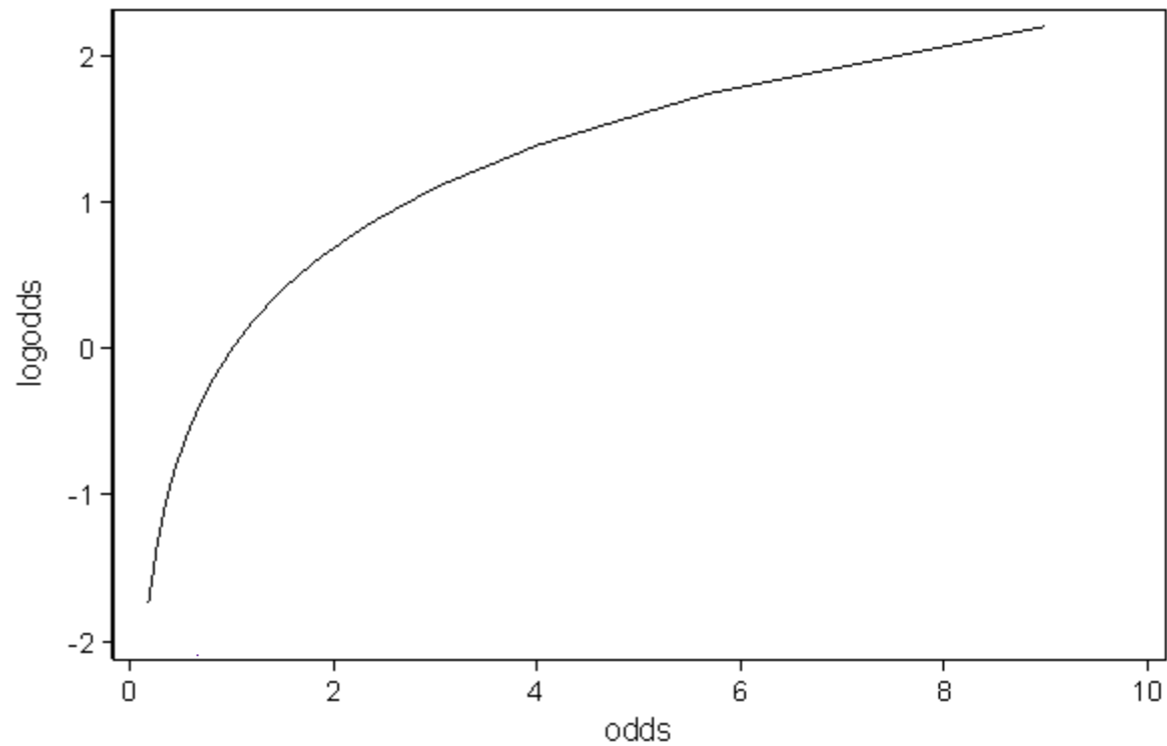
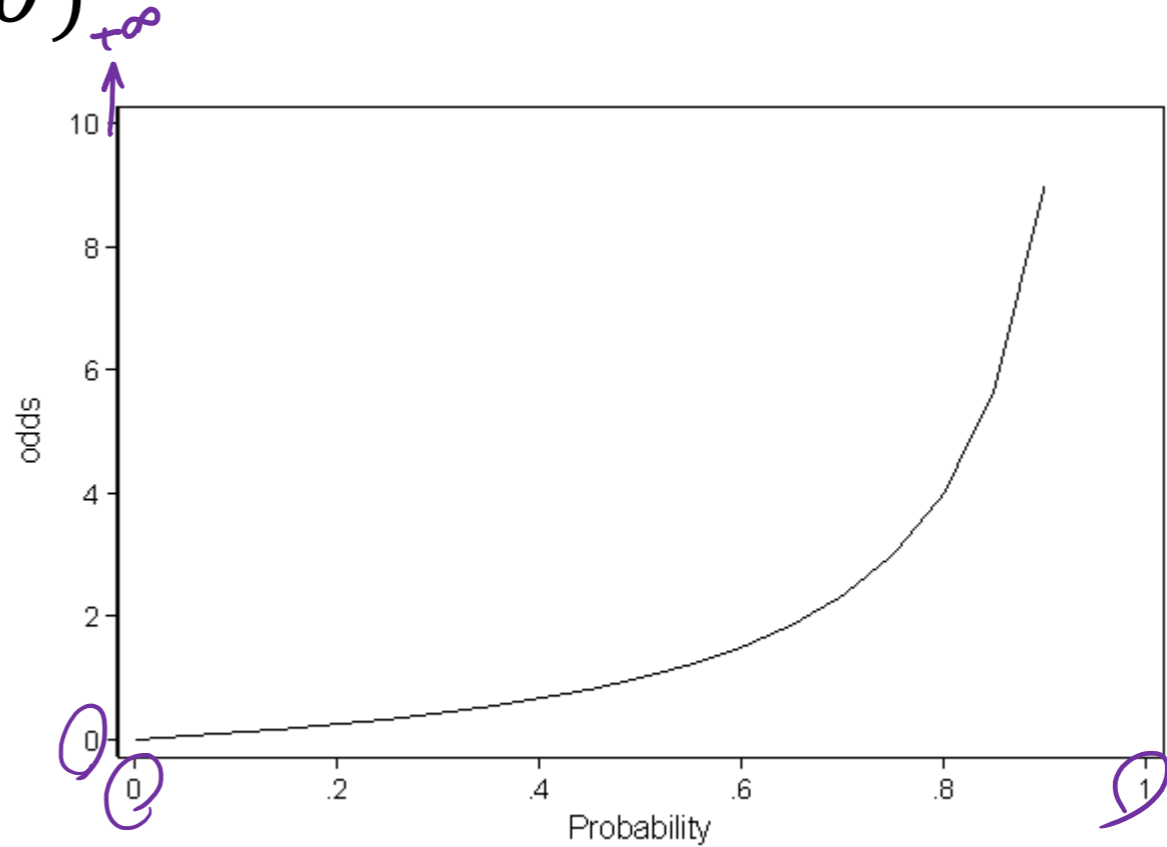
Summation

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



$\log(\text{odds})$ vs odds

What could be $x\theta$ domain?



Logistic function for posterior probability

$$-\infty < S = LCF = X\theta = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d < +\infty$$

$$g(s) = \frac{1}{1 + e^{-x^{\text{test}} \theta}} \begin{matrix} > 0.5 & \text{cat} \\ < & \text{dog} \end{matrix}$$

Let's use the following function:

$$s = x\theta$$

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

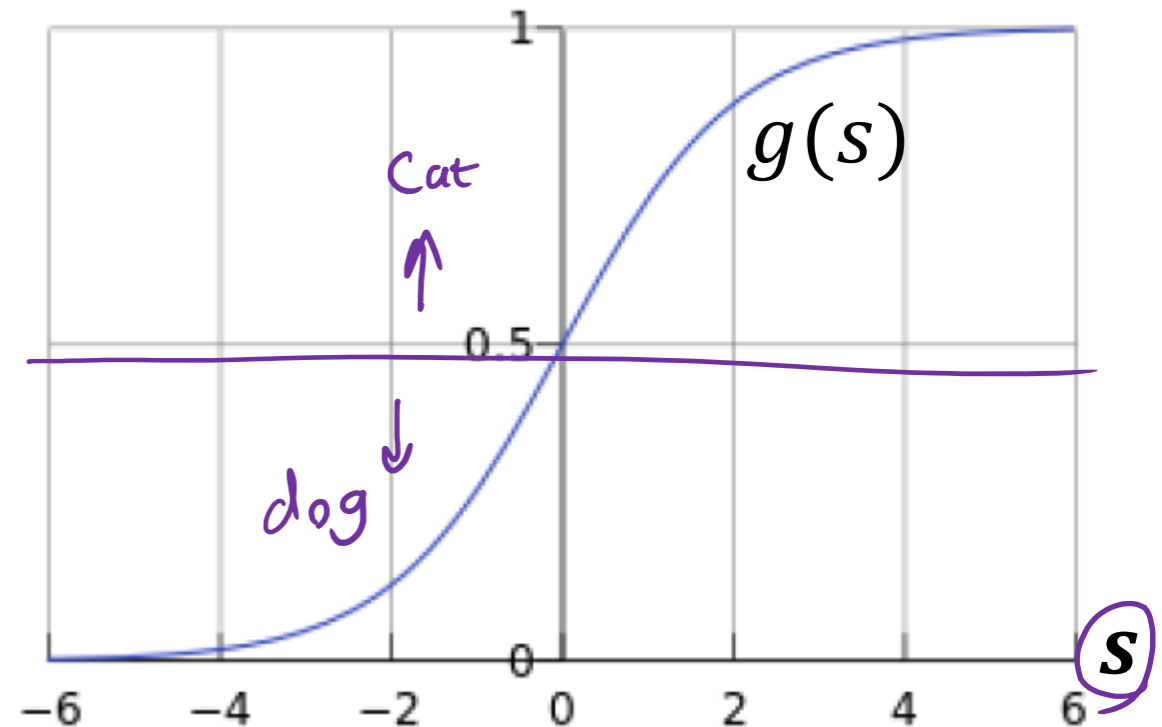
This formula is called sigmoid function

It is easier to use this function for optimization

Is 0.5 threshold cut-off a good choice?

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

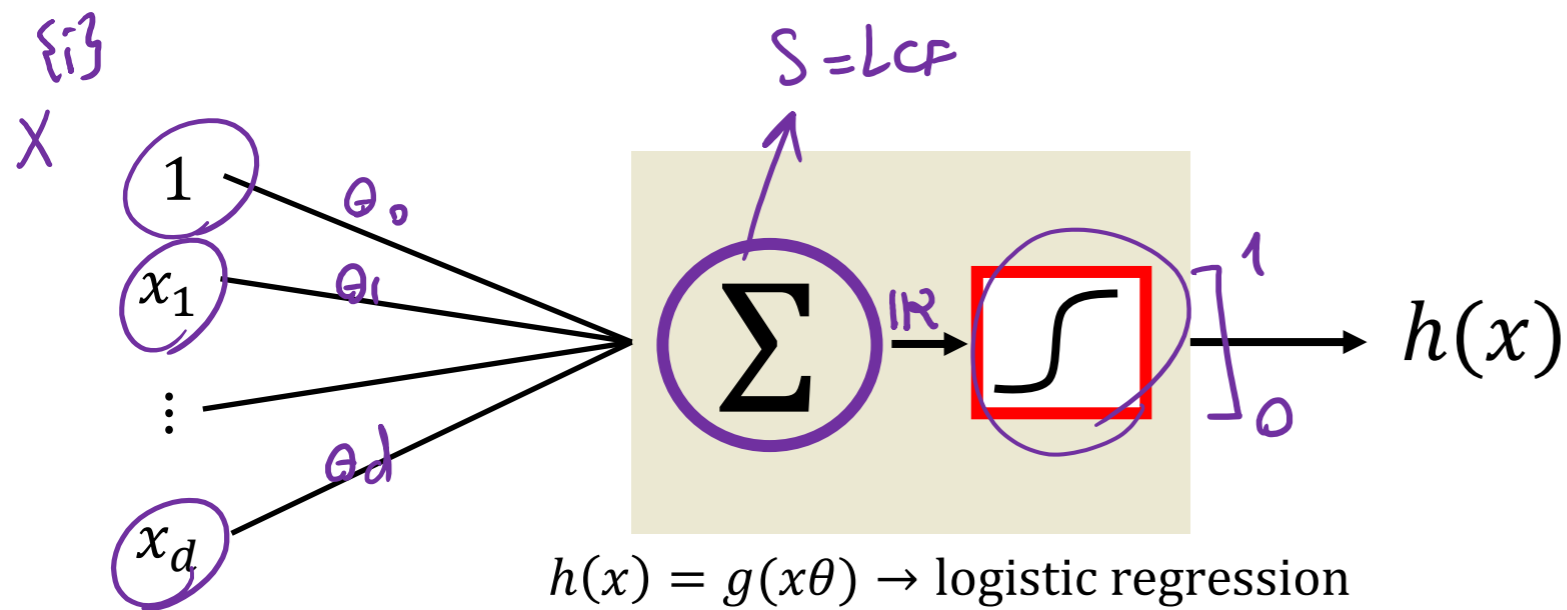
Many equations can give us this shape



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

Sigmoid Function

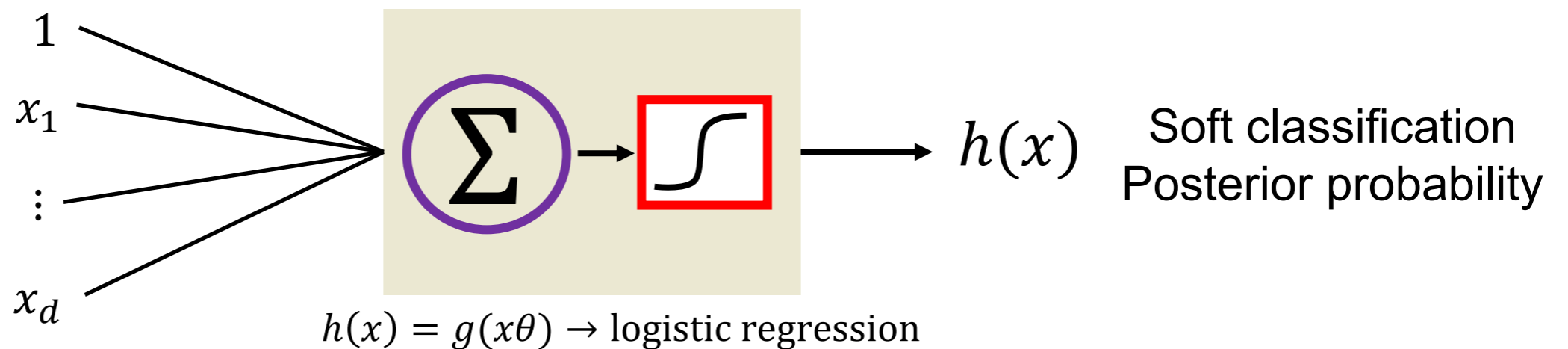
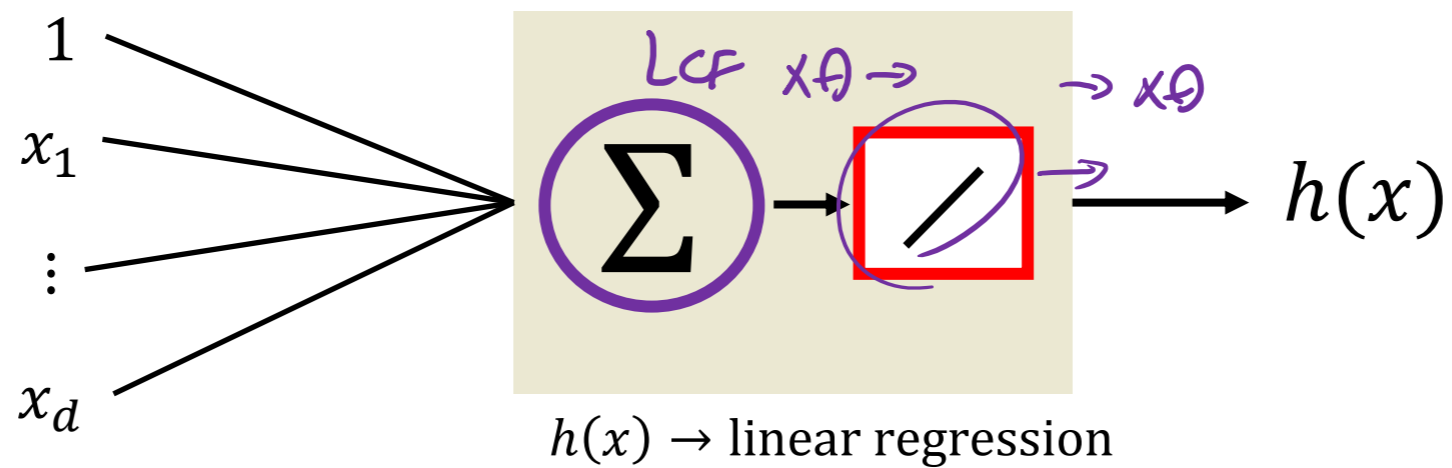
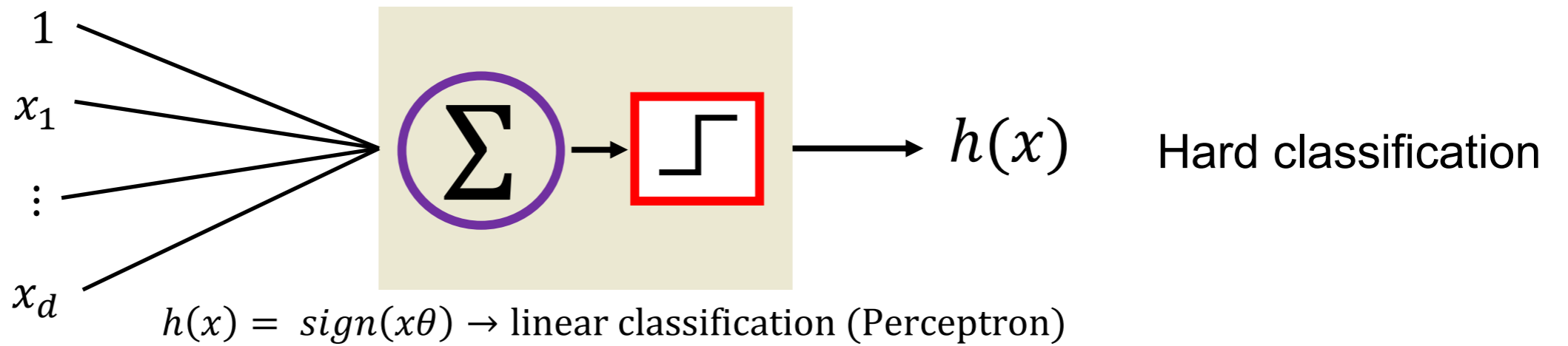
$$s = \sum_{i=0}^d x_i \theta_i = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$



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



$g(s)$ is interpreted as probability

Example: Prediction of heart attacks

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

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

We can't have a hard prediction here

$$s = x\theta$$

Let's call this risk score

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

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

Handwritten annotations:
- An arrow labeled "Sigmoid" points to $g(s)$.
- An arrow labeled "cat" points to the case $y = 1$.
- An arrow labeled "dog" points to the case $y = 0$.

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

Handwritten annotations:
- The term $p(y|x)$ is circled.
- The term $x\theta$ in the denominator is circled, with an arrow pointing to the label $d+1$.

Using posterior probability directly

Logistic regression model

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

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

Maximize \downarrow log likelihood

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

This form is concave, negative of this form is convex

The gradient of $l(\theta)$

$$\min_{\theta} l(\theta) = -\log \prod_{i=1}^n p(y^{i|} | x^{i|}, \theta) = CE = \sum y_a^{i|} \log y_p^{i|}$$

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

- Gradient

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

$$\Theta^{t+1} \leftarrow \Theta^t + \alpha \frac{\partial L}{\partial \theta}$$

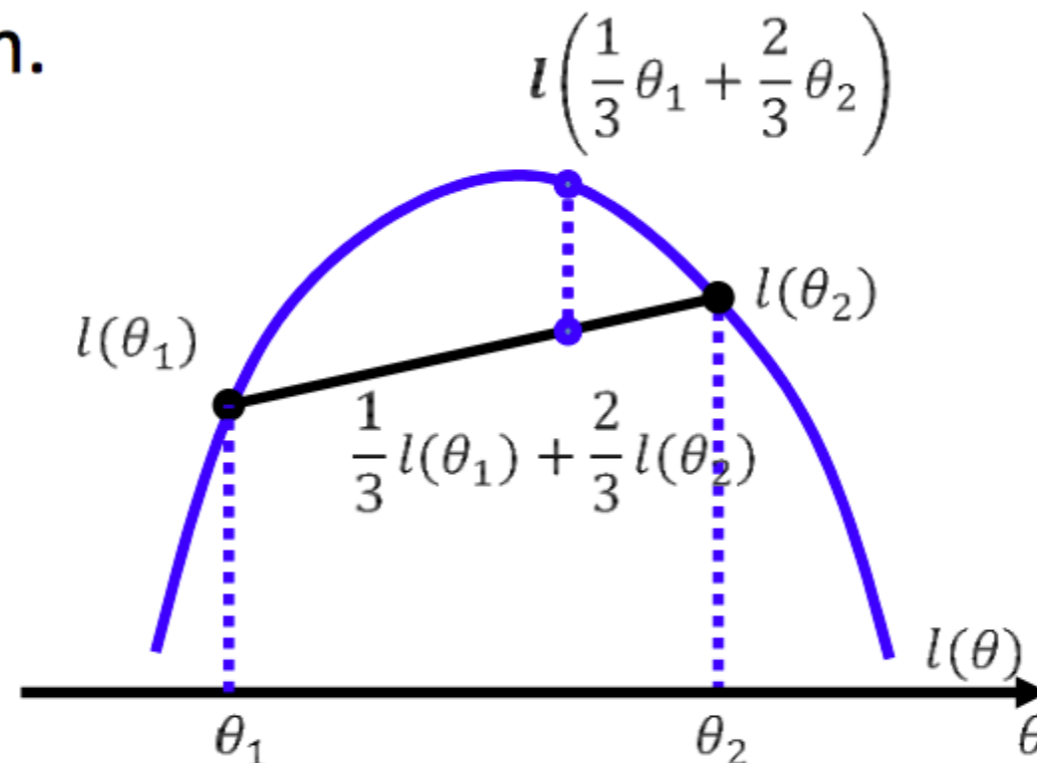
$\Theta^{t+1} \leftarrow \Theta^t - \alpha \frac{\partial L}{\partial \theta}$
 • Setting it to 0 does not lead to closed form solution

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.



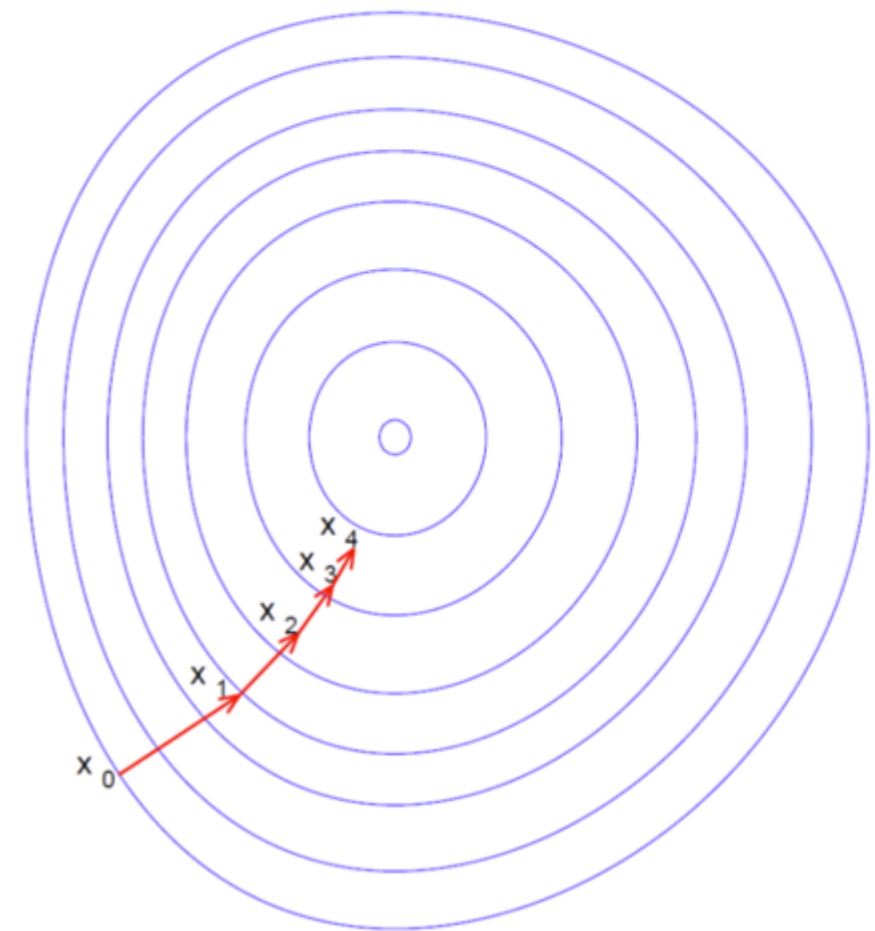
- Bad news: no closed form solution (resort to numerical method)

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

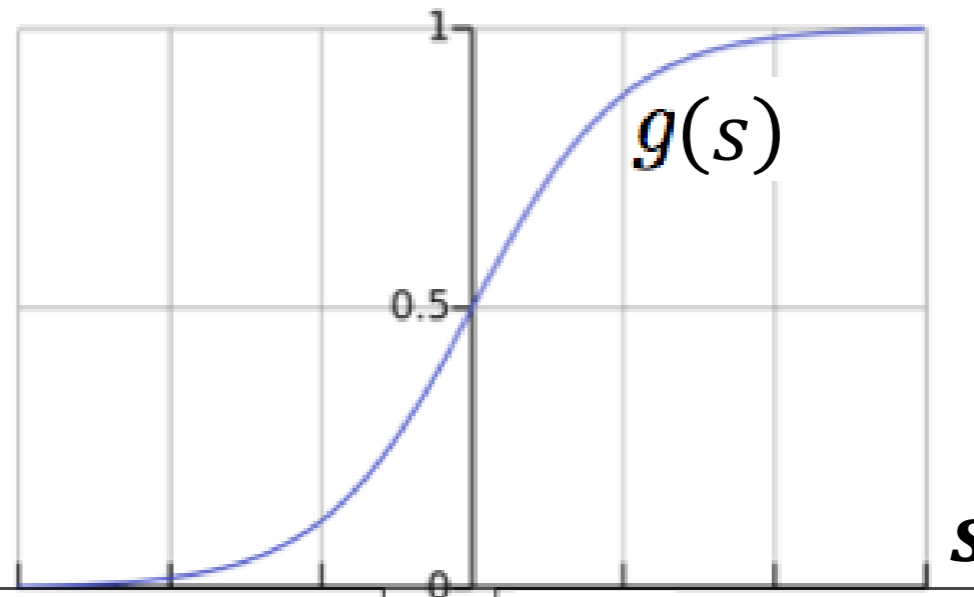
$$x_{k+1} = x_k - \gamma_k \nabla f(x_k)$$

γ_k is called the step size or learning rate



Logistic Regression

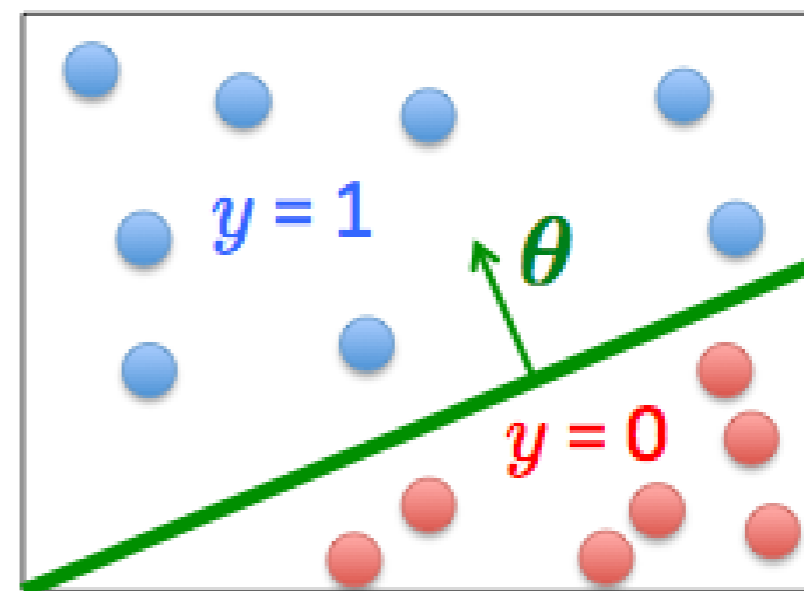
$$g(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$
$$s = x\theta$$



$x\theta$ 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 $g(s) < 0.5$

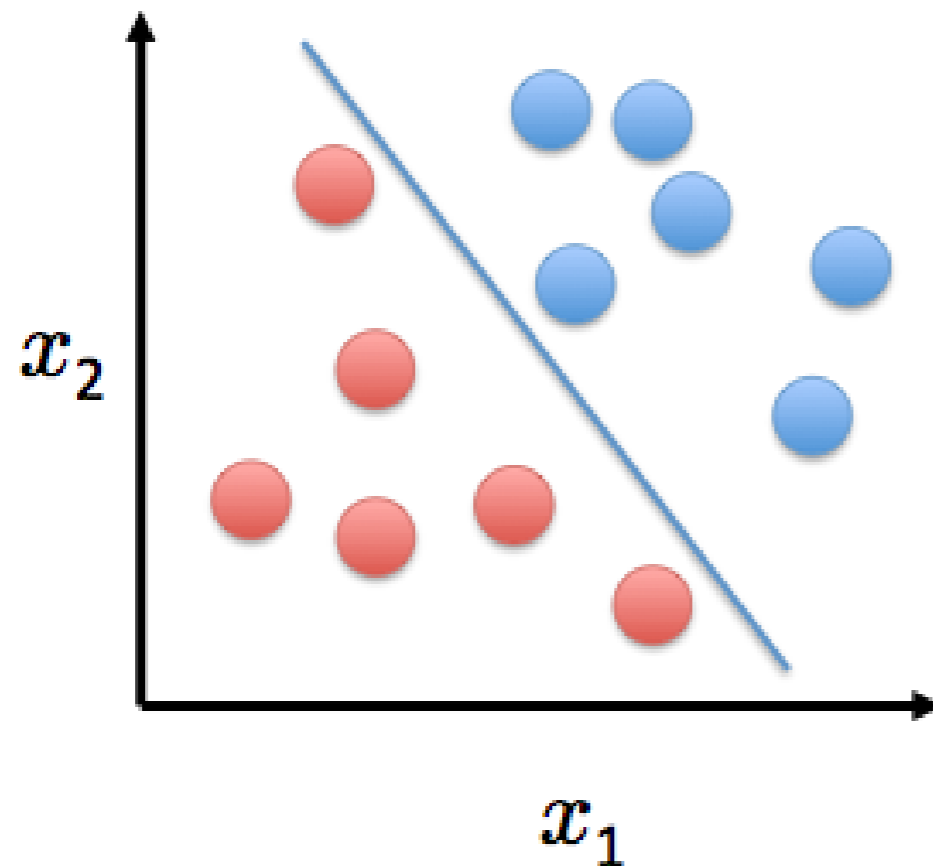


Outline

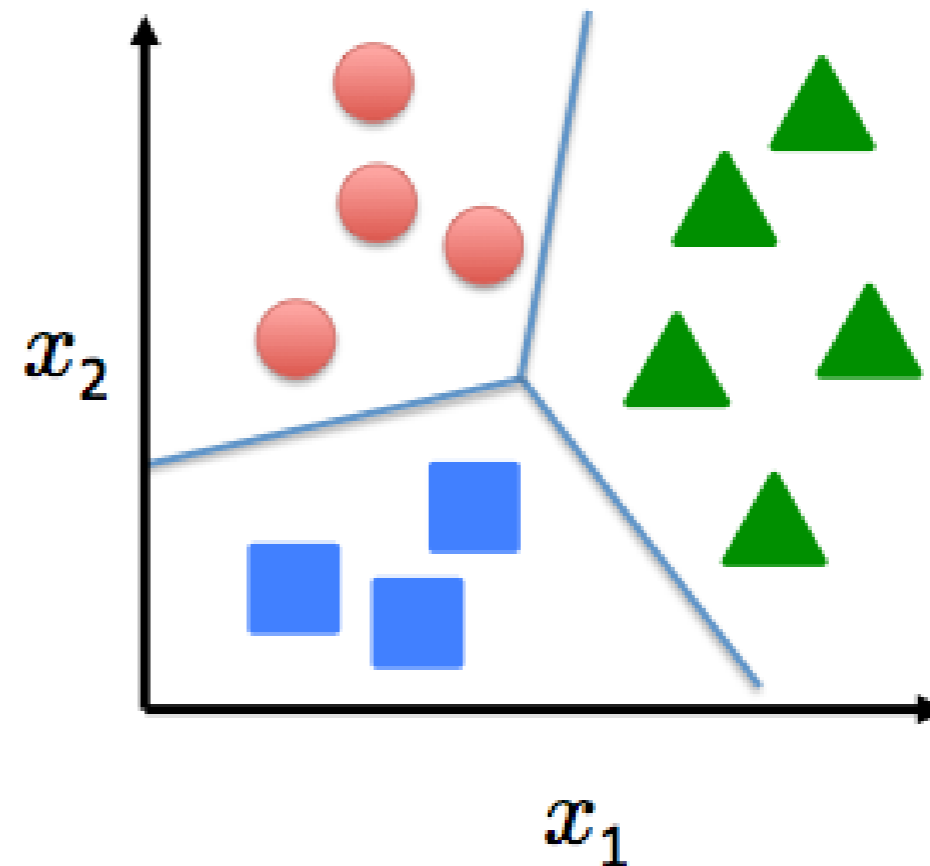
- Generative and Discriminative Classification
- The Logistic Regression Model
- Understanding the Objective Function
- Gradient Descent for Parameter Learning
- Multiclass Logistic Regression ←

Multiclass Logistic Regression

Binary classification:



Multi-class classification:



Disease diagnosis: healthy / cold / flu / pneumonia

Object classification: desk / chair / monitor / bookcase

Gradient Ascent(concave)/Descent(convex) algorithm

Negative log likelihood = Cross Entropy = $-\sum y_a^{i,j} \log y_p^{i,j}$

- Initialize parameter θ^0

$$\frac{\partial CE}{\partial \theta} = \bar{x}^T (y_p - y_a)$$

- Do

$$\frac{\partial L}{\partial \theta} \quad y_p = \frac{1}{1 + \exp(-s)} = \frac{1}{1 + \exp(-x\theta)}$$

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

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

$$\theta^t - \sum x^{i,j,T} (y_p^{i,j} - y_a^{i,j})$$

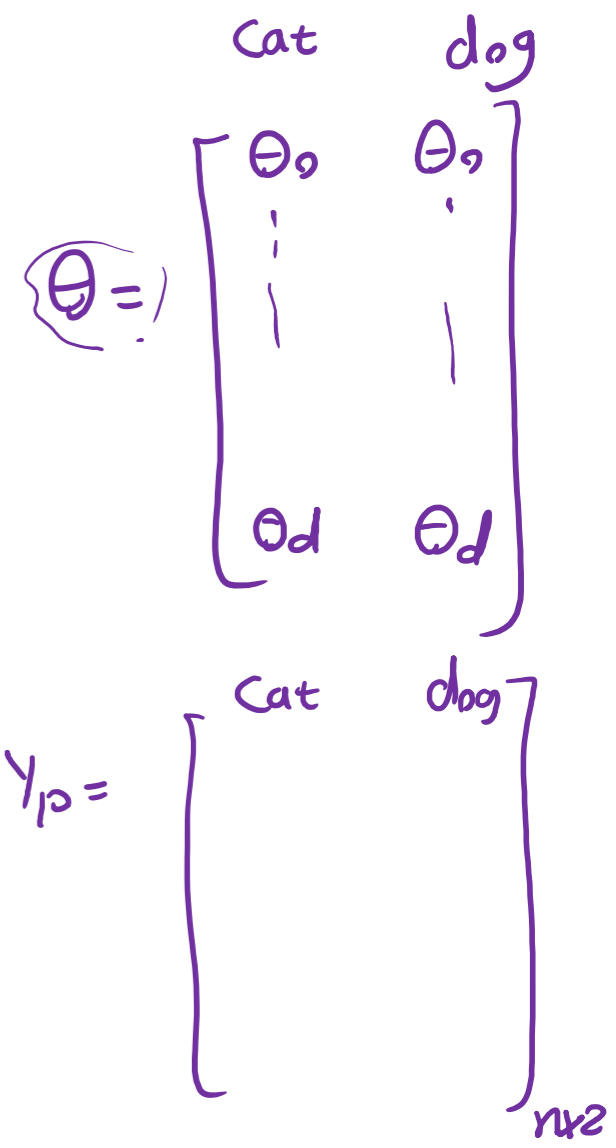
$$P(Y|X) = \text{Sigmoid} = \frac{1}{1 + \exp(-S)} \quad S = X\Theta$$

$$\log \prod_{i=1}^N P(y^{(i)} | x^{(i)})$$

$$P(Y=1|X) = \frac{\exp(S_1)}{\exp(S_1) + \exp(S_2)} \quad \text{cat}$$

$$S_1 = X\Theta_{\text{cat}} \text{ for cats}$$

$$S_2 = X\Theta_{\text{dog}} \text{ for dogs}$$



$$P(Y=2|X) = \frac{\exp(S_2)}{\exp(S_1) + \exp(S_2)} \quad \text{dog}$$

$$\Theta \leftarrow \Theta - \alpha X^T (Y_p - Y_a)$$

$d \times n$ \leftarrow $d \times n$ \leftarrow $n \times n$ \leftarrow $n \times 2$ \leftarrow $n \times 2$

$$P(Y=1|X) = \frac{1}{1 + \exp(S_2 - S_1)}$$

$$P(Y=2|X) = \frac{\exp(S_2 - S_1)}{1 + \exp(S_2 - S_1)}$$

$$S_2 - S_1 = X(\Theta_{\text{dog}} - \Theta_{\text{cat}}) = X\Theta$$

$$\Theta \leftarrow \Theta - \alpha X^T (Y_p - Y_a)$$

$$S_1 = X\theta_{cat} \quad S_2 = X\theta_{dog} \quad S_3 = X\theta_{fish}$$

$$Y_a = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \quad \hat{Y}_p = \begin{bmatrix} 0.7 & 0.2 & 0.1 \end{bmatrix}$$

dog
fish

$$\Theta = \begin{bmatrix} \theta_{cat} & \theta_{dog} & \theta_{fish} \end{bmatrix}$$

cat Reference

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

$$\frac{\exp(S_1)}{\exp(S_1) + \exp(S_2) + \exp(S_3)} = \frac{1}{1 + \underbrace{\exp(S_2 - S_1)}_{\theta_{(1)}} + \underbrace{\exp(S_3 - S_1)}_{\theta_{(2)}}}$$

dog

$$P(Y=2|x) =$$

fish

$$P(Y=3|x) =$$

$$\frac{\exp(S_2)}{C} = \frac{\exp(S_2 - S_1)}{1 + \exp(S_2 - S_1) + \exp(S_3 - S_1)} \quad (1)$$

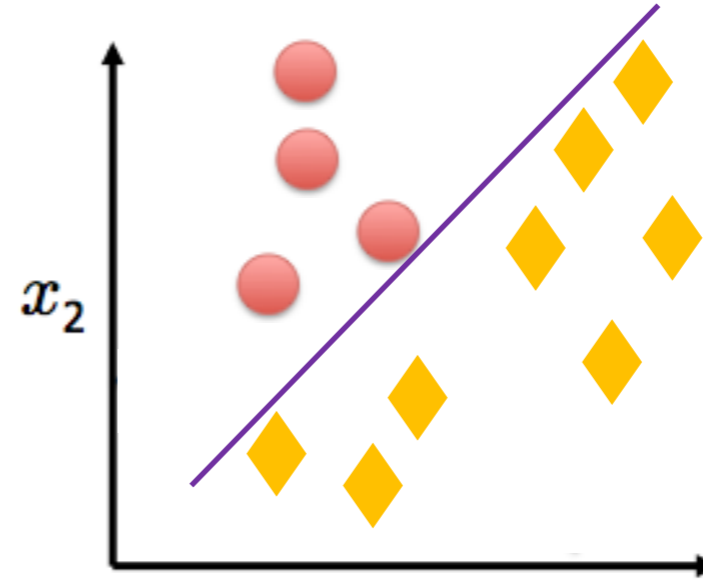
$$\frac{\exp(S_3)}{C} = \frac{\exp(S_3 - S_1)}{1 + \exp(S_2 - S_1) + \exp(S_3 - S_1)} \quad (2)$$

$$\Theta^{t+1} \leftarrow \Theta^t - \alpha \overline{X}^T (\hat{Y}_p - Y_a)$$

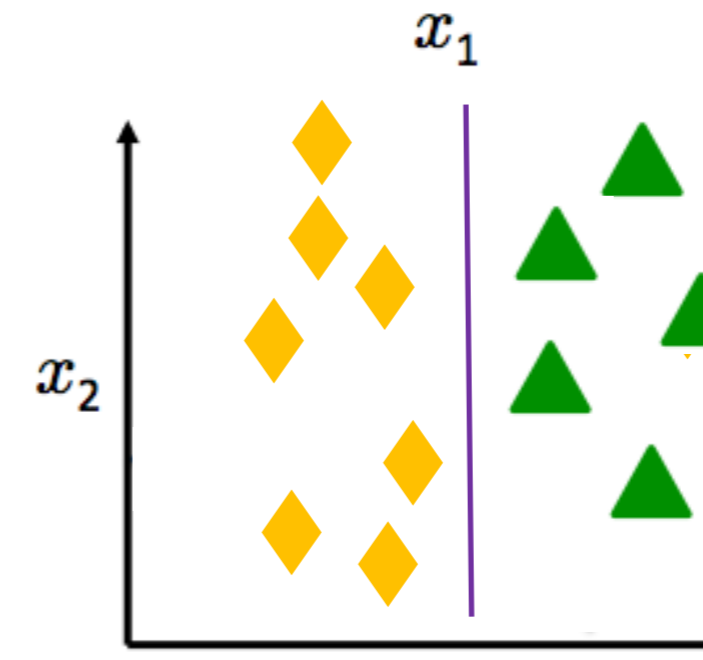
dog
fish

$N \times 2$
 $N \times 2$

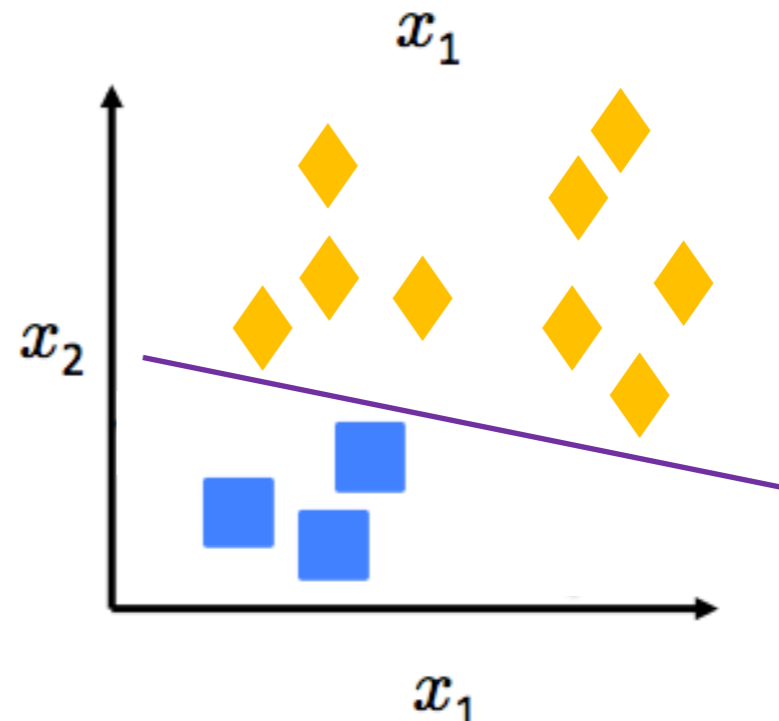
One-vs-all (one-vs-rest)



$h_{\theta}^1(x)$
0.6

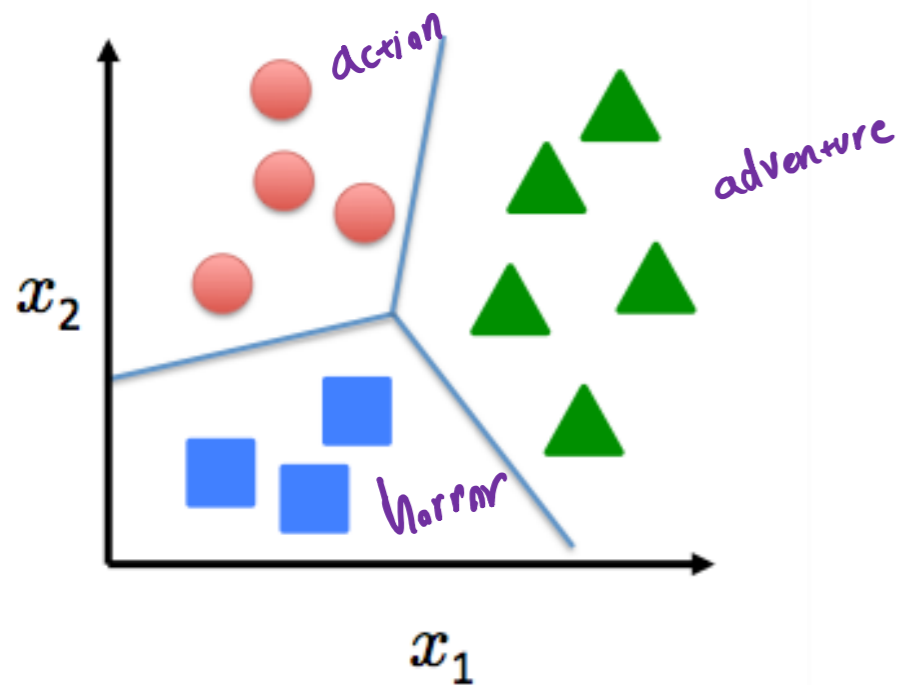


$h_{\theta}^2(x)$
0.75



$h_{\theta}^3(x)$
0.4

Multi-class classification:



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

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:

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

Using Softmax

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

$$y_a = [cat, dog, fish] = [1,0,0]$$

\Rightarrow there are M classes ($M = 3$ in this example)

$$y_p \text{ for class } m = \text{softmax}(x\theta) = \frac{\exp(x\theta)_m}{\sum_{j=0}^M \exp(x\theta)_j}$$

$$y_p = [0.6,0.3,0.1]$$

$$SGD \Rightarrow \theta^{t+1} \leftarrow \theta^t - \alpha \nabla L(\theta)$$

$$\theta^{t+1} \leftarrow \theta^t - \alpha x^T (y_p - y_a)$$

Take-Home Messages

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