

CS 229

Lecture 1 - Introduction

Lecture 2 - Linear Regression and Gradient Descent

• hypothesis $h(x) = \sum_{j=0}^n \theta_j x_j$

where $x_0 = 1$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

θ = "parameters"

m = # training examples

x = features

y = target $\Rightarrow (x, y)$: training example

$(x^{(i)}, y^{(i)})$ = i^{th} training example

n = # features

→ how to choose θ

$$h(x) \approx y \approx h_{\theta}(x)$$

$$\text{minimize } \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 =: J(\theta)$$

• Gradient Descent

Start w/ $\theta = \vec{\theta}$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \xrightarrow{\text{learning rate}}$$

$$\alpha := \alpha + 1$$

$$\rightarrow \frac{\partial}{\partial \theta_j} J(\theta) = (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y)$$

$$= (h_{\theta}(x) - y) \cdot x_j$$

$$\Rightarrow \theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

Repeat until convergence
(for $j = 1 \dots n$)

- Batch gradient descent → Compute whole Training Sample
↓ Slower but accurate
- Stochastic gradient → much faster
↓ Compute a single Training Sample

• Normal equation $\theta = R^{-1} b$

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix}$$

$$A \in \mathbb{R}^{m \times n}$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \cdots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$$

• example: $A \in \mathbb{R}^{2 \times 2}$

$$\text{tr } A = \text{trace of } A = \sum_i^n A_{ii}$$

$$f(A) = A_{11} + A_{22}^2$$

$$\Rightarrow \text{tr } A = \text{tr } A^T$$

$$\Rightarrow \nabla_A f(A) = \begin{bmatrix} 1 & 2A_{12} \\ 0 & 0 \end{bmatrix}$$

$$\nabla \rightarrow \nabla_{ij} = \delta_{ji}$$

but $i = j$

$$f(A) = \text{tr } AB \quad \text{fixed Matrix}$$

$$\nabla_A \text{tr } A A^T C = CA + C^T A$$

$$\nabla_A f(A) = B^T$$

$$\nabla_A \frac{1}{a} a^T c = Lac$$

$$\text{tr } A B = \text{tr } B A$$

$$\text{tr } ABC = \text{tr } CAB$$

Review these!

- elementwise operations $X\theta$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2} (X\theta - y)^T (X\theta - y) \quad \text{Error terms}$$

$$z^T z = \sum_i z_i^2$$

$$\nabla_{\theta} f(\theta) = \nabla_{\theta} \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

$$= X^T X \theta - X^T y \quad \stackrel{\text{Set } \rightarrow}{=} 0$$

$$\Rightarrow \theta = (X^T X)^{-1} X^T y \quad (\text{Normal function})$$

\hookrightarrow if X is non-invertible \rightarrow redundant feature

Lecture 3: locally weighted & logistic Regression

Recap:

$(x^{(i)}, y^{(i)})$ - i^{th} example

$x^{(i)} \in \mathbb{R}^{n+1}$, $y^{(i)} \in \mathbb{R}$, $x_0 = 1$

$m = \# \text{examples}$, $n = \# \text{features}$

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j = \theta^T x$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

\rightarrow locally weighted regression

"parametric" learning algorithm

\hookrightarrow fit fixed set of θ_i to data

"Non-parametric" learning algorithm

- ↳ Amount of data / θ s you need to keep grows (linearly) with size of data

→ LR: fit θ to minimize

$$\frac{1}{2} \sum_i^m (y^{(i)} - \theta^T x^{(i)})^2$$

Return $\theta^T x$

→ locally weighted regression: fit θ to minimize

$$\sum_{i=1}^m w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

where $w^{(i)}$ is a weighting function

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - z)^2}{2T^2}\right)$$

if $|x^{(i)} - z|$ is small, $w^{(i)} \approx 1$

$(x^{(i)}, y^{(i)})$ if $|z^{(i)} - z|$ is large, $w^{(i)} \approx 0$

↳ how much you should pay attention to $x^{(i)}$

T = "bandwidth"

→ Abbreviation

i.e. it is

e.g. for example

• Probabilistic interpretation

Why least squares?

→ Assume $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$

"error": unmodelled effects,
random noise

$$\epsilon^{(i)} \sim N(0, \sigma^2) \\ p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

→ this implies "parameterized by" not random $\overset{\text{mean}}{\uparrow}$

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

i.e. $y^{(i)} | x^{(i)}, \sigma \sim N(\theta^T x^{(i)}, \sigma^2)$ → variance

$$L(\theta) = p(y^{(i)} | x^{(i)}; \theta) \rightarrow \text{prob: } \theta \text{ fixed} \\ \text{vary data}$$

$$\text{"likelihood of } \theta \text{"} = \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta)$$

training fixed

$$\text{vary } \theta = \prod_{i=1}^n \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

[likelihood of parameters theta
[probability of data

$$\text{"log likelihood"} \\ l(\theta) = \log L(\theta) \uparrow \text{const} \\ = m \log \frac{1}{\sqrt{2\pi} \sigma} + \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$$

MLE: maximum likelihood estimator

Gaussian
IID

Choose θ to maximize $L(\theta)$

$$\text{i.e. choose } \theta \text{ to minimize: } \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 = J(\theta)$$

Classification problem

$y \in \{0,1\}$ (binary classification)

→ Logistic Regression

Want $h_\theta(x) \in [0,1]$ "sigmoid" or "logistic" function

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

$$\{ p(y=1|x; \theta) = h_\theta(x)$$

$$\{ p(y=0|x; \theta) = 1 - h_\theta(x)$$

$$\hookrightarrow P(y|x; \theta) = h(x)^y (1 - h(x))^{1-y}$$

$$\mathcal{L}(\theta) = P(\vec{y} | x; \theta)$$

$$= \prod_{i=1}^m h(x^{(i)})^{y^{(i)}} (1 - h(x^{(i)}))^{1-y^{(i)}}$$

$$\mathcal{L}(\theta) = \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1 - h_\theta(x^{(i)}))$$

→ Choose θ to $\overline{\max} \mathcal{L}(\theta)$

Batch gradient descent

$$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} \mathcal{L}(\theta)$$

$$= \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) \cdot x_j^{(i)}$$

$$= -\mathcal{L} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

Newton's Method

have f , want to find θ

$$\text{s.t.: } f(\theta) = 0 \quad (\text{i.e. want } f'(\theta) = 0)$$

$$\theta^{(t+1)} := \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})} \quad \text{"Quadratic convergence"}$$

(let $f(\theta) = 0$)

$$\theta^{(t+1)} := \theta^{(t)} - \frac{f'(\theta^{(t)})}{f''(\theta^{(t)})}$$

• When θ is a vector: $\theta \in \mathbb{R}^{n+1}$

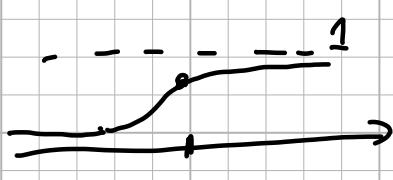
$$\theta^{(t+1)} := \theta^{(t)} + H^{-1} \underbrace{\nabla_{\theta} f}_{\rightarrow \text{vector } \mathbb{R}^{n+1}}$$

where H is the Hessian matrix

$$H_{ij} = \frac{\partial^2 f}{\partial \theta_i \partial \theta_j} \rightarrow \text{expensive}$$

Lecture 4 Perceptron & Generalized Linear Model

logistic Regression sigmoid



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

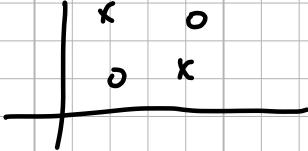
$$g(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}$$

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

$$\hat{\theta}_j := \theta_j + \lambda (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

(not used in practice) λ

e.g.



• Exponential family

$$\text{PDF} : p(y; \eta) = b(y) \exp \left[\eta^T T(y) - a(\eta) \right]$$

y - data

η - eta : natural parameter

$T(y)$ $\stackrel{?}{=} y$ sufficient statistic ?

$b(y)$ - base measure

$a(\eta)$ - log partition

$$\frac{\phi^y \cdot (1-\phi)^{1-y}}{\lambda(1-\phi)^y}$$

• Bernoulli (Binary Data)

$$= \exp \left(\log \left(\phi^y (1-\phi)^{1-y} \right) \right)$$

ϕ : prob of event

$$p(y; \phi) = \phi^y (1-\phi)^{1-y} = \exp \left(\log \left(\frac{\phi}{1-\phi} \right) \cdot y + \log(1-\phi) \right)$$

$$\begin{cases} b(y) = 1 \\ n = \log \left(\frac{\phi}{1-\phi} \right) \end{cases} \Rightarrow \phi = \frac{1}{1 + e^{-n}}$$

$$\begin{cases} T(y) = y \\ a(\eta) = -\log(1-\phi) \end{cases} \Rightarrow -\log \left(1 - \frac{1}{1 + e^{-n}} \right) = \log \left(1 + e^{-n} \right)$$

• Gaussian (with fixed variance)

Assume $\sigma^2 = 1$

$$\begin{aligned} P(y; \mu) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right) \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} \underbrace{\left(\mu - \frac{1}{2} \mu^2\right)}_{\eta} \sim a(\eta) \\ &\quad b(y) \end{aligned}$$

Properties:

@ MLE wrt $\eta \Rightarrow$ Convex

NLL is convex

$$\textcircled{1} \quad f[y; \eta] = \frac{\partial}{\partial \eta} a(\eta) \quad \left. \begin{array}{l} \text{No} \\ \text{integrals} \end{array} \right\}$$

$$\textcircled{2} \quad \text{Var}[y; \eta] = \frac{\partial^2}{\partial \eta^2} a(\eta)$$

• GLM

Assumption / Design choices

(i) $y | x; \theta \sim \text{Exponential Family } (\eta)$

(ii) $\eta = \theta^T x \quad \theta \in \mathbb{R}^n, x \in \mathbb{R}^w$

(iii) Test time: output $E[y | x; \theta] = h_\theta(x)$
 $= E[y; \theta^T x]$

mean is the predict

$$\hookrightarrow \max_{\theta} \log p(y^{(j)}; \theta^T \eta^{(j)})$$

GLM training

→ learning update rule

$$\theta_j := \theta_j + \alpha (y^{(j)} - h_{\theta}(x^{(j)})) x_j^{(j)}$$

Terminology

η : natural parameter

$E[y; \eta] = g(\eta) \rightarrow$ Canonical response function

$\eta = g^{-1}(u) \rightarrow$ Canonical link function

$$g'(\eta) = \frac{\partial}{\partial \eta} a(\eta)$$

• 3 - parameterizations

Model Param	Natural Param	Canonical Param
θ	η	ϕ - Ber
\uparrow learn - $\theta^T \eta$	$\xrightarrow{g^{-1}}$	$\mu \sigma^2$: Gauss
	$\leftarrow \frac{1}{g}$	λ : Poisson

Design choice

$$\bullet \text{Logistic Regression: } h_{\theta}(x) \cdot E(y|x; \theta) = \phi = \frac{1}{1 + e^{-\eta}} = \frac{1}{1 + e^{-\theta^T x}}$$

Real - Gaussian

Binary - Bernoulli

Count - Poisson

R^+ - Gamma, Exponential

Distn - Beta, Dirichlet

Bayesian

- Gaussian $\eta = \mu$
 - Softmax Regression - Cross Entropy
- Multiclassification

k : # classes label $y \in \{0,1\}^k$ one-hot
vector
 $x^{(i)} \in \mathbb{R}^n$

$$\hat{p}(y) - p(y)$$

$$\begin{aligned}\text{Cross Ent} (p, \hat{p}) &= - \sum_{y \in \{0,1,\dots\}} p(y) \cdot \log \hat{p}(y) \\ &= - \log \hat{p}(y_0) \\ &= - \log \frac{e^{\alpha_0^T x}}{\sum_{c \in \{0,1,\dots\}} e^{\alpha_c^T x}}\end{aligned}$$

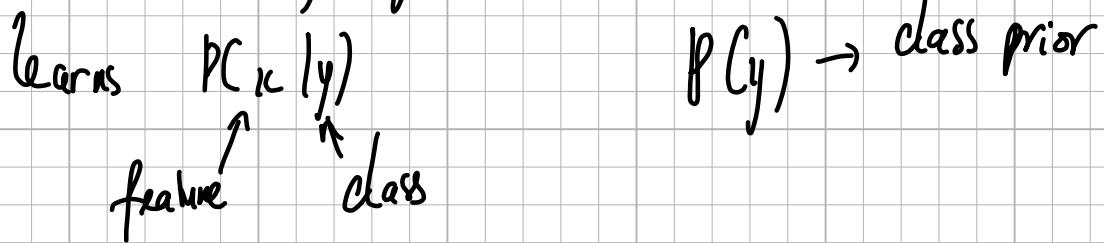
Lecture 5 GDA & Naive Bayes

- Discriminative Learning Algorithms
- Gaussian Discriminant Analysis

- Generative { Discriminative composition
- Naive Bayes

→ Discriminative: learn $p(y|x)$
 (or learn $h_0(x) = \langle \cdot, 1 \rangle$ directly)

→ Generative learning algo.



• Bayes rule:

$$P(y=1|x) = \frac{P(x|y=1) \cdot p(y=1)}{P(x)}$$

$$P(x) = P(x|y=1)p(y=1) + P(x|y=0)p(y=0)$$

• Gaussian Discriminatory Analysis (GDA)

Suppose: $x \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

Assume $p(x|y)$ is Gaussian

d : dimension of matrix

[multivariate Gaussian $f(\mathbf{x}) =$

uni-variate $\frac{1}{(2\pi)^{d/2} |\Sigma|^{\frac{1}{2}}} \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$

$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma) \quad \mathbf{z} \in \mathbb{R}^n$$

\downarrow \downarrow
 \mathbb{R}^n $\mathbb{R}^{n \times n}$

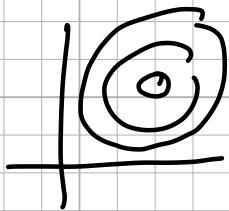
$$\rightarrow E[\mathbf{z}] = \boldsymbol{\mu}$$

$$\rightarrow \text{Cov}(\mathbf{z}) = E[(\mathbf{z} - \boldsymbol{\mu})(\mathbf{z} - \boldsymbol{\mu})^\top]$$

$$= E[\mathbf{z}\mathbf{z}^\top] - (E[\mathbf{z}]) (E[\mathbf{z}])^\top$$

$$\rightarrow p(\mathbf{z}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \cdot (\mathbf{z} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu})\right)$$

l.g. $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ Contour of Gaussian density



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- GDA model. Parameters: $\mu_0, \mu_1, \Sigma, \phi$

$$P(x|y=0), P(x|y=1) \quad R^w \quad R^{n \times n} \quad R$$

$$P(y) = \phi^y (1-\phi)^{1-y}$$

$$\rightarrow \text{Training set } \left\{ (x^{(i)}, y^{(i)}) \right\}_{i=1}^n$$

Joint likelihood:

$$L(\phi, \mu_0, \mu_1, \Sigma) = \prod_{i=1}^n P(x^{(i)}|y^{(i)}). P(y^{(i)})$$

Discrimination:

(Maximum likelihood)

$$\mathcal{L}(\theta) = \prod_{i=1}^n P(y^{(i)}|x^{(i)}; \theta)$$

Maximize likelihood estimation

$$\phi = \frac{\sum_{i=1}^n y^{(i)}}{m} = \frac{\sum_{i=1}^m 1 \{ y^{(i)} = 1 \}}{m}$$

$$\mu_0 = \frac{\sum_{i=1}^m 1 \{ y^{(i)} = 0 | x^{(i)} \}}{m}$$

$$\sum_{i=1}^m 1 \{ y^{(i)} = 0 \}$$

$Tg \mu_0 \mu_1$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_y^{(i)}) (\mu_x^{(i)} - \mu_y^{(i)})^T$$

→ Prediction $\left[\min_z (z - \xi)^2 = 0 \quad \arg\max_z (z - \xi)^2 = \xi \right]$

$$\arg \max_y p(y|x) = \arg \max_y \frac{p(x|y) p(y)}{p(x)}$$

$$= \arg \max_y p(x|y) p(y)$$

Generative \Rightarrow Discriminative
 Gaussian logistic Regression
 (stranger) (weaker)

→ Categorical Test and gain conviction
 or just use logistic Regression

→ Naive Bayes

feature vector x ? $\lambda = \begin{bmatrix} \lambda_1 & \dots & \lambda_n \end{bmatrix} \quad \lambda_i \in \{0, 1\}^n$
 $\lambda_i = 1 \text{ if word } i \text{ appears in mail}$
 2^{16500} possible values

Assume $x_i | \Sigma$ are conditional independent
 $P(x|y) = \prod_{i=1}^m P(x_i|y)$

o Parameters $q_j|y=1 = P(x_j = 1 | y = 1)$

$$q_j|y=0 = P(x_j = 1 | y = 0)$$

$$\phi_y = p(y=1)$$

Joint likelihood:

$$L(\phi_y, \phi_j|y) = \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi_y, \phi_j|y)$$

MLE:

$$\hat{\phi}_y = \frac{1}{n} \sum_{i=1}^n \zeta_y^{(i)}$$

$$\hat{\phi}_j|y=1 = \frac{\sum_{i=1}^n \zeta_j^{(i)}}{\sum_{i=1}^n \zeta_j^{(i)}}$$