

1/19/2023: MLE for Linear Regression, Logistic Regression

$$L(w) = \frac{1}{n} \sum_{i=1}^n (w^T x^{(i)} - y^{(i)})^2$$

← why not 4?
← why not absolute value

Maximum Likelihood Estimation

- Post a probabilistic process that generated our data
- Find parameters that make observed data seem most likely

Coin Flips (testbed for MLE)

[Heads, tails, heads, heads, heads] ← observed data

p = probability of flipping heads once ← unknown parameter

eg. if $p = \frac{1}{3}$: $\frac{1}{3} \times \frac{2}{3} \times \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3}$ ← likelihood given value of p

Linear Regression: Assume $y^{(i)}$ drawn from

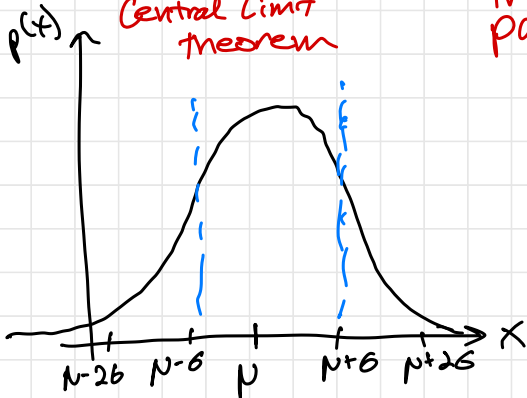
Gaussian w/ mean $w^T x^{(i)}$ & variance σ^2

Central Limit theorem

true parameter

independently

constant



$$p(x; \mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$P(A|B)$

Likelihood of data

$$\begin{aligned}
 \mathcal{L}(w) &= \prod_{i=1}^n P(y^{(i)} | x^{(i)}; w) \\
 &= \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right)
 \end{aligned}$$

"parameterized by"

maximize $\mathcal{L}(w)$ is equivalent to maximizing $\log \mathcal{L}(w)$

$$\log \mathcal{L}(w) = \sum_{i=1}^n \left[\log \frac{1}{\sigma \sqrt{2\pi}} + \left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2} \right) \right]$$

log-likelihood

constant

$$= \underbrace{-\frac{1}{2\sigma^2}}_{\text{neg}} \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2$$

maximizing $\log \mathcal{L}(w)$ is equivalent to minimizing $L(w)$

$$L(w) = \underbrace{\frac{1}{n}}_{\text{pos}} \sum_{i=1}^n (w^T x^{(i)} - y^{(i)})^2$$

linear regression loss

Classification

Goal: predict "label"/"class" from a discrete set of options

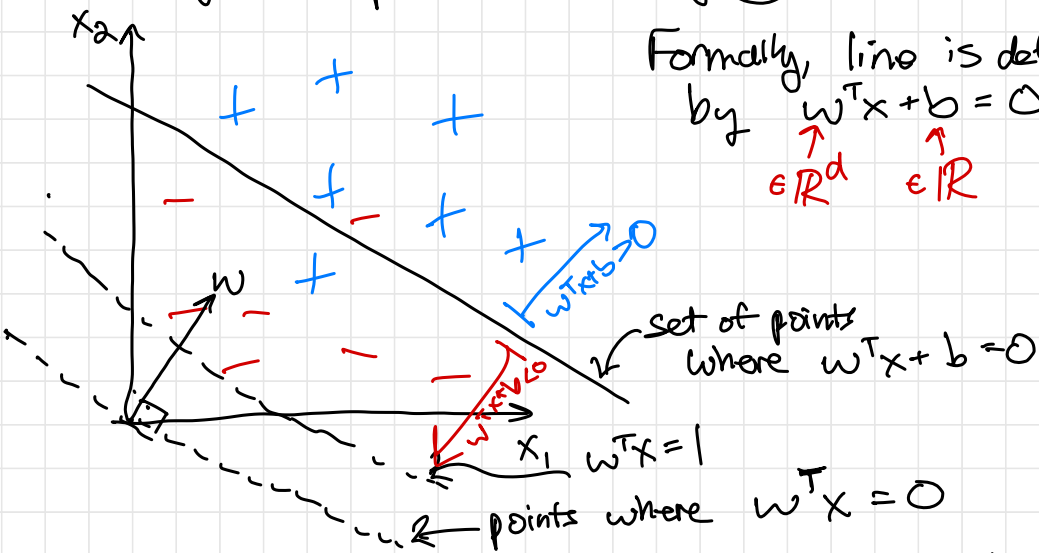
- Tumor: Benign or Malignant
 - Email: spam vs not spam
 - Handwritten digits: is it 0, 1, 2, ..., 9
- Binary classification (2 labels)
Multi-class classification (> 2 labels)

Binary classification terms:

one label is "positive", other is "negative"

$y = 1$ $y = -1$ (or $y = 0$)

Modelling Assumption of linearity (⊖)



Formally, line is defined by $w^T x + b = 0$

$w \in \mathbb{R}^d$ $b \in \mathbb{R}$

Predictions: IF $w^T x + b > 0$, predict $y = 1$
if $w^T x + b < 0$ predict $y = -1$

Use MLE to come up with appropriate loss function
(We can omit for same reason as before)

$$p(y=1 | x; w) = \frac{1}{1 + \exp(-w^T x)} = \sigma(w^T x)$$

↑ "sigmoid" or "logistic" func.

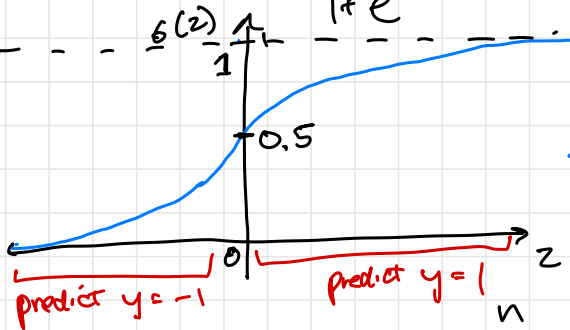
where $\sigma(z) = \frac{1}{1 + e^{-z}}$

Convenient fact:

$$p(y|x; w) = \sigma(y w^T x)$$

If $y=1$ true

If $y=-1$: $\sigma(-w^T x)$
 $= \frac{1}{1 + e^{w^T x}} = \frac{\exp(-w^T x)}{\exp(-w^T x) + 1}$



Maximize $\mathcal{Q}(w) = \log \prod_{i=1}^n p(y^{(i)} | x^{(i)}; w)$
 $= \sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; w)$
 $= \sum_{i=1}^n \log \sigma(y^{(i)} w^T x^{(i)})$

equivalent to minimizing

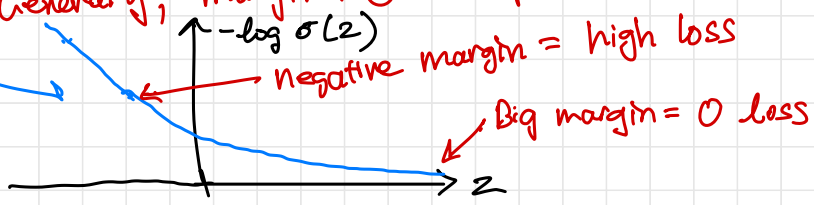
$$J(w) = \sum_{i=1}^n -\log \sigma(y^{(i)} w^T x^{(i)})$$

* the loss function for logistic regression

- If $y^{(i)} = 1$, want $w^T x^{(i)}$ large
- If $y^{(i)} = -1$, want $w^T x^{(i)}$ small

"margin"
large margin is good

- Generally, margin $> 0 \iff$ prediction is correct



Fact:
 $-\log \sigma(z)$
 is
 convex

Gradient descent time

$J(w)$ is convex (by same rules as linear reg.)

$$J(w) = \sum_{i=1}^n \underbrace{-\log \sigma}_{\text{positive number}} (y^{(i)} w^T x^{(i)})$$

$$\frac{d}{dz} -\log \sigma(z) = -\sigma(-z)$$

$$\nabla_w J(w) = \sum_{i=1}^n \underbrace{-\sigma(-y^{(i)} w^T x^{(i)})}_{\text{positive number}} \times \underbrace{y^{(i)}}_{\text{+/- ?}} \cdot \underbrace{x^{(i)}}_{\text{vector}}$$

If $y^{(i)} = 1$, adding multiple of $x^{(i)}$ to w
makes $w^T x^{(i)}$ larger
increases $p(y^{(i)} | x^{(i)}; w)$

If $y^{(i)} = -1$, subtract multiple of $x^{(i)}$ from w

$\sigma(-\text{margin})$

If margin large: ≈ 0 ← Already doing well

If margin small: ≈ 1 ↑ Room for improvement

} Prioritization