

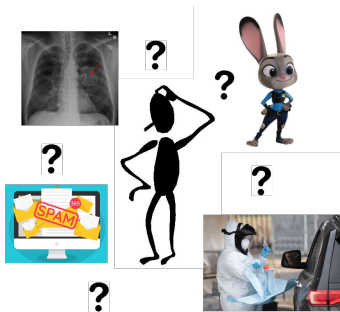
2. Logistic regression

- Binary classification
- Logistic regression
- This is a broad overview, we will revisit each section in depth

Binary classification

Mind-reading task

- Will Alice like the movie?
- Will Bob click on this link?
- Is this cancer?
- Is this spam?
- Is this a terrorist?



Features: (age, height, political leanings, historical decisions, friends)

Labels: (yes or no)

Mind-reading task

Features: (age, height, political leanings, historical decisions, friends)

Labels: (yes or no)

Mind-reading task

Features: (age, height, political leanings, **historical decisions**, friends)

Labels: (yes or no)

- Alice liked Finding Nemo, so she'll like Zootopia

Mind-reading task

Features: (age, height, **political leanings**, historical decisions, friends)

Labels: (yes or no)

- Alice liked Finding Nemo, so she'll like Zootopia
- Bob and Alice are both Republican. Bob likes this tweet → so will Alice

Mind-reading task

Features: (age, **height**, political leanings, historical decisions, friends)

Labels: (yes or no)

- Alice liked Finding Nemo, so she'll like Zootopia
- Bob and Alice are both Republican. Bob likes this tweet → so will Alice
- Claire and Dennis are both tall. Claire likes to ski → so will Dennis.

Mind-reading task

Features: (age, height, political leanings, historical decisions, friends)

Labels: (yes or no)

- Alice liked Finding Nemo, so she'll like Zootopia
- Bob and Alice are both Republican. Bob likes this tweet → so will Alice
- Claire and Dennis are both tall. Claire likes to ski → so will Dennis.

First approach: linear model

$$\underbrace{\text{label}}_y = \text{sign} \left(\sum_{k=1}^d \overset{\text{theta}}{\text{weight}_k} \times \overset{x}{\text{feature}_k} \right)$$

Training: learn weights so that prediction on training set is about right

Logistic regression model

Logistic “regression” model

- Not regression (wrongly named) but classification
- Features $x[1] = \text{age}$, $x[2] = \text{height}$, $x[3] = \text{past viewing behavior...}$
- We write $x \in \mathbb{R}^d$ to mean x is a vector with d values
- Importance weights $\theta \in \mathbb{R}^d$

june july aug.
↓ ↓ ↓
[1, 1, 0]

$$\theta[i] = \begin{cases} \text{large, positive} & = \text{feature correlates with label} \\ \text{small, around 0} & = \text{feature probably not that important} \\ \text{large, negative} & = \text{feature inversely correlates with label} \end{cases}$$

Logit model

$$\text{Log likelihood} = \log \left(\frac{\Pr(y = 1|x, \theta)}{\Pr(y = 0|x, \theta)} \right) = \underbrace{\theta^T x := \sum_{k=1}^d \theta[k]x[k]}_{\text{notation for inner product}}$$

$\theta \cdot x$

example of correlated 1-D data + labels

$x_1 = .1, x_2 = 1., x_3 = 2., x_4 = -5.$

$y_1 = 1, y_2 = 1, y_3 = 1, y_4 = 0$

Logistic regression model

- Logit model

$$\text{Log likelihood} := \log \left(\frac{\Pr(y = 1|x, \theta)}{\Pr(y = 0|x, \theta)} \right) = x^T \theta$$

- Rearrange, normalize $\Pr(y = 1) + \Pr(y = 0) = 1$

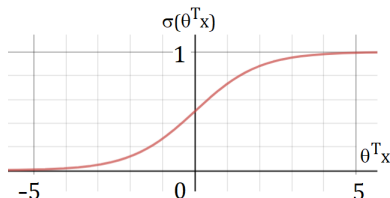
$$\Pr(y = 1|x, \theta) = \frac{e^{\theta^T x}}{1 + e^{\theta^T x}} =: \sigma(\theta^T x)$$

- $\sigma(\theta^T x)$ is the sigmoidal function, models “soft probability”

[cat, dog, camel, bird]

1, 0, 0 0

0, 1, 0 0



Discriminative vs generative

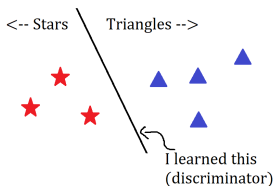
Logistic regression is a discriminative approach

$$\Pr(y = 1|x, \theta) = \sigma(\theta^T x), \quad \Pr(x|y, \theta) = ??? \text{ (don't care)}$$

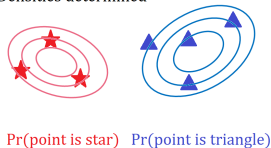
Discriminative

vs

generative



Densities determined



Discriminative vs generative

Logistic regression is a discriminative approach

$$\Pr(y = 1|x, \theta) = \sigma(\theta^T x), \quad \Pr(x|y, \theta) = ??? \text{ (don't care)}$$

Discriminative

vs

generative



...roses are red,



violets are blue...

Know the symptoms of COVID-19, which can include the following:



Gauge by symptoms severity



"Deep dream"

@dreamscapeapp on Twitter

This bird is black with green and has a very short beak



Zhang et al. StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks

Discriminative vs generative

Logistic regression is a discriminative approach

$$\Pr(y = 1|x, \theta) = \sigma(\theta^T x), \quad \Pr(x|y, \theta) = ??? \text{ (don't care)}$$

Discriminative	vs	generative
<ul style="list-style-type: none">-- if that's all the task requires (classification only)--may have less parameters<ul style="list-style-type: none">** don't need as much data		<ul style="list-style-type: none">--can be used to create more training data<ul style="list-style-type: none">** depends on distribution being good-- more detailed/complete feature representation-- easy to screw up by outliers

Logistic regression

- Training samples: $\mathcal{X} = \{x_1, \dots, x_m\} \subset \mathbb{R}^d$ features, $\mathcal{Y} = \{y_1, \dots, y_m\}$ labels
- yes : $y_k = 1$, no : $y_k = 0$
- We assume the data is independently, identically distributed (i.i.d.)
- Logit model: $\Pr(y = 1|x) = \sigma(\theta^T x)$
- Goal: find maximum likelihood estimator

likelihood = $\Pr(\text{Thing you see} \mid \text{underlying factor})$

$$\theta_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \underbrace{\Pr(\mathcal{X}, \mathcal{Y} \mid \theta)}_{\Pr(Y \mid X, \theta)}$$

- Predictor for new sample x

$$\tilde{y}_{\text{pred}} = \mathbf{sign}(x^T \theta_{\text{MLE}})$$

$\{-1, 1\}$

$$\begin{aligned} \text{tilde } y &= 2 * y - 1 \\ y &= \text{tilde } y / 2 + .5 \end{aligned}$$

Derive logistic regression objective function

Likelihood of label

$$\begin{aligned} \Pr(y = 1 \mid x, \theta) &= \sigma(x^T \theta) \\ \Pr(y = 0 \mid x, \theta) &= 1 - \sigma(x^T \theta) \end{aligned}$$

$$\begin{aligned} \Pr(\mathcal{Y} \mid \mathcal{X}, \theta) &\stackrel{\text{i.i.d.}}{=} \prod_{i=1}^m \Pr(y_i \mid x_i, \theta) \quad y_i \in \{0, 1\} \\ &\stackrel{\text{Logit model}}{=} \prod_{i=1}^m \sigma(\theta^T x_i)^{y_i} (1 - \sigma(\theta^T x_i))^{1-y_i} \end{aligned}$$

Maximum log likelihood

$$\begin{aligned} \max_{\theta} \log(\Pr(\mathcal{Y} \mid \mathcal{X}, \theta)) &\iff \max_{\theta} \sum_{i=1}^m \underbrace{y_i \log \sigma(\theta^T x_i)}_{\text{nonzero if } y_i=1} \\ &\quad + \underbrace{(1 - y_i) \log(1 - \sigma(\theta^T x_i))}_{\text{nonzero if } y_i=0}, \quad y_i \in \{0, 1\} \end{aligned}$$

$$\sigma(s) = e^s / (1 + e^s)$$

$$\begin{aligned} \sigma(-s) &= 1 - \sigma(s) \\ &\iff \end{aligned}$$

convex, $0 \leq a \leq 1$

$f(a * x + (1-a) * y) \leq a * f(x) + (1-a) * f(y)$

f is concave if $-f$ is convex

$$\max_{\theta} \sum_{i=1}^m \log \sigma(\tilde{y}_i x_i^T \theta), \quad \tilde{y}_i = 2y_i - 1 \in \{-1, 1\}$$

$$x_i^T \theta + 1 * \theta[0]$$

$$x_i = [1, \dots]$$

Maximum likelihood estimator

$$y \sigma(x.T \theta) = \sigma(y x.T \theta)$$

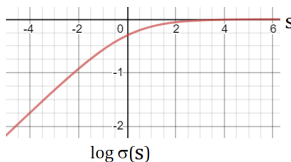
$$\theta_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \left(\mathcal{L}(\theta; \mathcal{X}, \mathcal{Y}) := \sum_{i=1}^m \log \sigma(y_i x_i^T \theta) \right), \quad y_i \in \{-1, 1\}$$

rewards

$y_i * \sum_{k=1}^d x_i[k] * \theta[k]$

margin

really good margin



- function $\mathcal{L}(\theta; \mathcal{X}, \mathcal{Y})$ is concave in θ

- θ_{MLE} are the stationary points of \mathcal{L}

$$\theta = \theta_{\text{MLE}} \iff \nabla_{\theta} \mathcal{L}(\theta; \mathcal{X}, \mathcal{Y}) = 0$$

if and only if



Training

$$\theta = \theta_{\text{MLE}} \iff \nabla_{\theta} \mathcal{L}(\theta; \mathcal{X}, \mathcal{Y}) = 0$$

Gradient ascent

- For a function $\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}$, the gradient $\nabla \mathcal{L}(\theta) \in \mathbb{R}^d$,

$$\nabla \mathcal{L}(\theta) = \left(\frac{\partial \mathcal{L}(\theta)}{\partial \theta_1}, \frac{\partial \mathcal{L}(\theta)}{\partial \theta_2}, \dots, \frac{\partial \mathcal{L}(\theta)}{\partial \theta_d} \right)$$

- Keep climbing up! $L(t) = 1 + t_1 + t_2^2$, $g L(t) = (1, 2 t_2)$

$$\theta^{(0)} = \text{anywhere}, \quad \theta^{(k+1)} = \theta^{(k)} + \alpha \nabla \mathcal{L}(\theta^{(k)})$$

for a suitably small step size $\alpha > 0$

- Stop when $\nabla \mathcal{L}(\theta^{(k)})$ is close enough to 0

test misclass rate (k) - test misclass rate(k-1) \approx 0

Extensions: acceleration, using higher order derivatives, parallelization, stochastic gradients...

Logistic regression summary

- **Training:** Using data x_i , labels $y_i \in \{-1, 1\}$, find

$$\theta = \arg \min_{\theta} \sum_{i=1}^m \log \sigma(y_i x_i^T \theta)$$

- **Prediction:** New data sample x

$$y = \begin{cases} \text{sign}(\theta^T x) & \text{(hard decision)} \\ \sigma(\theta^T x) & \text{(soft decision)} \end{cases}$$

Generalization

- Predictor θ_{MLE} depends crucially on training set \mathcal{X}, \mathcal{Y}
- But what if training sample is not that representative?
 - Not enough data (and coverage for rare events)
 - Presence of damaging outliers
 - Data corrupted, anonymized, or tampered
- Given loss function $\mathcal{L}(\theta; x, y)$ promoting high $\Pr(y|\theta, x)$, finite training set $\mathcal{T} = \{(x_1, y_1), \dots\}$

$$\begin{aligned} R^* &= \underbrace{\mathbb{E}_{x,y}[\mathcal{L}(\theta; x, y)]}_{\text{Expected risk}}, & R_{\mathcal{T}} &= \underbrace{\frac{1}{|\mathcal{T}|} \sum_{(x,y) \in \mathcal{T}} \mathcal{L}(\theta; x_i, y_i)}_{\text{Empirical risk, training loss}} \\ R^* &= \underbrace{R_{\mathcal{T}}}_{\text{get}} + \underbrace{(R^* - R_{\mathcal{T}})}_{\text{Generalization loss}} \end{aligned}$$

training loss

want

- Solutions?: regularization, MAP estimator, ensemble learning, more data ...

What else? We will cover

- Further analysis on logistic regression
 - multiclass classification, margin maximization method, ...
- Other classification methods
 - thresholded linear regression, support vector machines, decision trees ...
- More details on computation of training
 - how to pick step size, stochastic methods, nonconvex models ...
- Generative approaches
 - GMMs, HMMs, expectation maximization ...
- ...

Important but we will sweep under the rug

- Data balancing (rare diseases, natural disasters, car accidents)
- Data preprocessing/cleaning
- Data anonymizing/privacy
- Cost of making the wrong decision (ethical, computational)
- Other competing metrics (cost vs quality of service, retention vs addiction)

