

Xiao and Suya

task T performance measure P experience E
 (data)

linear regression:

task - predict y output, given X input
 $\hat{y} = w^T x$

Performance measure (P) - mean square error (on test data)
 MSE_{test}

experience (E): training data set $(X; Y)$
 Split into $(X, Y)_{train}$ $(X, Y)_{test}$
 Cost function minimize $MSE_{train} := \arg \min_w MSE_{train}(w)$

logistic regression: classification task - output is probability of class
 (binary)

$$(X, y): P(y=1 | x, \theta) = \sigma(\theta^T x)$$

Sigmoid

$$P(y=0 | x, \theta) = 1 - \sigma(\theta^T x)$$

typically, predict output class by max probability

decision boundary still linear -

$$\frac{1}{1+e^{-\theta^T x}} > 0.5$$

$$e^{-\theta^T x} < 1$$

$$\theta^T x > 0$$

training points $\{x^{(1)}, \dots, x^{(m)}\}$ This requires i.i.d. assumption

$$\text{likelihood } L(\theta) = \prod_{i=1}^m p(y_i | x_i, \theta)$$

$$= \prod_{i=1}^m (\sigma(\theta^T x_i))^{y_i} \cdot (1 - \sigma(\theta^T x_i))^{1-y_i}$$

MLE: $\arg \max_{\theta} L(\theta)$

use log likelihood: $\arg \min_{\theta} -\log(L(\theta))$ avoid underflow problems with product

why is MLE better than MSE?

- ① consistency no bias
- ② always can find global minimum (convex)



MSE is non-convex

gradient descent to minimize:

$$\begin{aligned} \theta_{ML} &= \arg \min_{\theta} -\log L(\theta) \\ &= \arg \min_{\theta} \sum_{i=1}^m -\log \left[\underbrace{-y_i \log(\sigma(\theta^T x_i))}_{J_i(\theta)} - (1-y_i) \log(1-\sigma(\theta^T x_i)) \right] \\ \frac{\partial J_i(\theta)}{\partial \theta} &= \begin{bmatrix} \frac{\partial J_i(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J_i(\theta)}{\partial \theta_n} \end{bmatrix} \end{aligned}$$

$\left[\frac{\partial J_i(\theta)}{\partial \theta_j} \right]$

using chain rule

$$\frac{\partial J_i(\theta)}{\partial \theta_j} = -y_i \frac{1}{\sigma(\theta^T x)} \cdot \frac{\partial (\sigma(\theta^T x))}{\partial \theta_j} - (1-y_i) \frac{-1}{1-\sigma(\theta^T x)} \cdot \left(\frac{-\partial (\sigma(\theta^T x))}{\partial \theta_j} \right)$$

$\sigma'(z) = \frac{e^{-z}}{(1+e^{-z})^2}$ need to compute gradient
derivative of $\sigma(z) = \frac{1}{1+e^{-z}}$

$$\begin{aligned} \frac{\partial (\sigma(\theta^T x))}{\partial \theta_j} &= \sigma(\theta^T x) (1-\sigma(\theta^T x)) \\ &= -y_i (1-\sigma(\theta^T x_i)) x_i - (1-y_i) \sigma(\theta^T x_i) (-x_i) \\ &= -y_i x_i + \sigma(\theta^T x_i) x_i \\ &= \sigma(\theta^T x_i - y_i) x_i \end{aligned}$$

so, $J(\theta) = -\log L(\theta)$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^n [\sigma(\theta^T x_i - y_i) x_i]$$

initialization $\theta = 0$
for $t = 1, \dots, T$ $\theta = \theta - \sum \frac{\partial J(\theta)}{\partial \theta_j}$ gradient descent
output θ

guarantees?

if convex, with low enough ϵ , always converges
if non-convex, no guarantees

regularization: add term to avoid overfitting

L^2 norm regularization $\lambda \|\theta\|_2^2$
 $\underset{\theta}{\operatorname{argmin}} J(\theta) + \lambda \|\theta\|_2^2$

or $\underset{\theta}{\operatorname{argmin}} J(\theta)$
subject to constraint $\|\theta\|_2 \leq C$ constrained optimization

rewrite as unconstrained:
 $\underset{\theta}{\operatorname{argmin}} J(\theta) + \lambda \|\theta\|_2^2 - \lambda C$ constant

L^1 vs L^2 regularization
 $\underset{\theta}{\operatorname{argmin}} f(\theta) + \lambda \|\theta\|_2^2 \rightarrow \underset{\theta}{\operatorname{argmin}} f(\theta)$ will constraint $\|\theta\|_2^2 \leq C$

$\underset{\theta}{\operatorname{argmin}} f(\theta) + \lambda \|\theta\|_1 \rightarrow \underset{\theta}{\operatorname{argmin}} f(\theta)$ with constraint $\|\theta\|_1 \leq C$.

L_1 norm constraint $x_1^2 + x_2^2 \leq C$

L_2 norm constraint $|x_1| + |x_2| \leq C$

L_1 not differentiable
constraint vs. unconstrained
 L_1 pushes to corner \rightarrow sparse (some 0s)
 L_2 preserves high dimensionality