

Regression overview, least-squares and logistic Regression

Classification: given pt. x , predict class

Regression: given pt. x , predict numerical value

- Choose regression fn $h(x; p)$ w/ parameters p (h = hypothesis)
- optimize a cost fn (based on lost fn)
 - ex. risk fn

fns:

① linear: $h(x; w, \alpha) = w^T x + \alpha$

② polynomial

③ logistic: $h(x; w, \alpha) = S(w^T x + \alpha)$; $S(r) = \frac{1}{1 + e^{-r}}$ logistic fn

↳ often used for estimating probabilities because its always between 0 and 1

↳ natural way for LDA (can get w/o LDA "logistic regression")

loss fns: z is prediction $h(x)$; y is true label

Ⓐ $L(z, y) = (z - y)^2$ (squared error)
↳ amplifies outliers

Ⓑ $L(z, y) = |z - y|$ (abs error)
↳ hard to optimize

③ $l(z, y) = -y \ln z - (1-y) \ln(1-z)$ (logistic loss)
 \hookrightarrow cross entropy; $y \in [0, 1]$, $z \in (0, 1)$

Cost fns:

① $J(h) = \left[\frac{1}{n} \right] \sum_{i=1}^n l(h(x_i), y_i)$ (mean loss)
 \uparrow doesn't affect optimization

② $J(h) = \max_{i=1}^n l(h(x_i), y_i)$ (max loss)
 \uparrow trustworthy data - no outliers to throw off

③ $J(h) = \sum_{i=1}^n w_i l(h(x_i), y_i)$ (weighted sum)
 \uparrow some data more important, so we weight more

④ $J(h) = (a, b, c) + \underbrace{\lambda \|w\|_2^2}_{\text{regularization term}} \rightarrow$
 $\rightarrow l_2$ penalized
 \rightarrow dec. in multi sol cases

⑤ " " " " $\lambda \|w\|_1 \rightarrow$
 $\rightarrow l_1$ penalized
 \rightarrow encourage sparsity

famous regression methods:

⇒ least squares linear regression: ① - ① - ②

$$\frac{1}{n} \sum_{i=1}^n \|W^T X_i + \alpha - y_i\|^2$$

⇒ " " but weighted: ① - ① - ③

$$\sum_{i=1}^n w_i \|W^T X_i + \alpha - y_i\|^2$$

⇒ Ridge Regression: ① - ① - ④

⇒ LASSO: ① - ① - ⑤] QP

⇒ logistic regression: ③ - ③ - ⑤] convex; min w/ grad descent

⇒ least abs. deviations: ① - ③ - ⑥] LP

⇒ Chebyshev criterion: ① - ③ - ⑥] LP

least squares linear regression

$$\min_{w, \alpha} \frac{1}{n} \sum_{i=1}^n (W^T X_i + \alpha - y_i)^2 \Rightarrow \arg \min_{w, \alpha} \sum_{i=1}^n \|W^T X_i + \alpha - y_i\|^2$$

X is $n \times d$ design matrix, y is n -vector of scalar labels

X_i^T is a point of dim d

X_{ij} is a feature column

after derivatives dim trick:
 X is $n \times d+1$; w $d+1$ vector

$$[X_1 \ X_2 \ 1] \begin{bmatrix} w_1 \\ w_2 \\ \alpha \end{bmatrix} \rightarrow h(x) = w^T X$$

$$\Rightarrow \arg \min_w \|Xw + y\|^2 \quad (\text{new opt prob, same as before})$$

RSS(w), residual sum of squares

Solving: $w^T X^T X w + 2 w^T X^T y + y^T y$

$$\nabla_w f = 2 X^T X w - 2 X^T y = 0 \Rightarrow X^T X w = X^T y$$

= normal equations

if $X^T X$ singular, problem unconstrained

if not, then $w^* = \underbrace{(X^T X)^{-1} X^T y}_{\text{linear transformation}}$

$X^+ X = I$ (left inv.) \hookrightarrow pseudo inverse X^+

$$\hat{y}_i = w^T X_i \Rightarrow \hat{y} = Xw = XX^+ y = Hy \quad H := \text{hat matrix}$$

Advantages:

- easy to compute
- Unique, stable solution

disadvantages:

- sensitive to outliers
- $X^T X$ should be singular

Logistic Regression

- logistic regression fn, fits probs (0,1)
- usually used for classification
- discriminative model: jump to fitting, skip gaussians

$$\arg \min_w \sum_{i=1}^n l(X_i^T w, y_i) = - \sum_{i=1}^n (y_i \ln s(X_i^T w) + (1-y_i) \ln (1-s(X_i^T w)))$$

$J(w)$ is convex! Solve by gradient descent

$$s(\gamma) = \frac{e^{-\gamma}}{(1+e^{-\gamma})^2} = s(\gamma)(1-s(\gamma))$$

$$\begin{aligned}
s_i &= S(X_i^T w) & \nabla_w J &= - \sum \left(\frac{y_i}{s_i} \nabla s_i - \frac{1-y_i}{1-s_i} \nabla s_i \right) \\
& & &= - \sum \left(\frac{y_i}{s_i} - \frac{1-y_i}{1-s_i} \right) s_i (1-s_i) X_i \\
& & &= - \sum (y_i - s_i) X_i \\
& & &= - X^T (y - S(Xw)) & S(Xw) &= \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{bmatrix}
\end{aligned}$$

gradient descent rule:

$$w \leftarrow w + \epsilon X^T (y - S(X^T w))$$

Stochastic:

$$w \leftarrow w + \epsilon (y_i - S(X_i^T w)) X_i$$

Converges before we visit all pts for large n

Start: $w=0$ works well